

Supplementary Information

Figure Captions

SI Figure 1. Orthogonal NMF produces low-error reconstruction Orthogonal nonnegative matrix factorization factors each a) gene expression data matrix into b) F , a set of orthogonal linear feature vectors, and c) C , a matrix of coefficients. For (b) we order the genes according to their weight in each feature, keeping only genes with greater $> 4\sigma$ for each feature. This ordering highlights the orthogonality of features. (d) The reconstructed data FC is a denoised version of the original data matrix. (e) Data values from the original data matrix are partitioned into bins, and their means are plotted against corresponding data values from the reconstructed data matrix (FC). Bin width ~ 0.1 . Shaded area indicates the standard deviation between original values and reconstructed values within each bin. (See Methods - oNMF error analysis). (f) Histogram of the entries in normalized data matrix D , with $\log(\text{counts})$ on y-axis. 99.996% of data values are less than 5, which is also the linear regime of (e). (g) Zoomed in view of white inset from (c), showing specific genes that have high weights for two features, as well as their top gene set results from gene set enrichment analysis.

SI Figure 2. Choosing dimensionality of feature set Plot of loss function for different feature sets, built by sweeping over many values of m (x-axis). Feature set with lowest loss (red asterisk) is chosen. This example uses data from Fig. 4. Loss function is a function of the residual error and a penalty on m . For loss function equation, see Methods - Extraction of gene feature vectors with matrix factorization.

SI Figure 3. PopAlign models for all 12 tissues of Tabula Muris For each tissue, we plot the following: (upper left) joint t-SNE plot of experimental single-cell data (black), PopAlign model-generated data (teal), and mixture model centroids (μ) as numbered disks, (upper right) Heatmap of the proportion of cells classified by each mixture model component against their annotated labels. Each column sums to 1. (lower) Heatmap of mixture component centroids (μ) in terms of their expression level of features. For all tissues, we used a universal 30D feature set determined using sampled data from all tissues within the collection.

SI Figure 4. Mixture components across Tabula Muris tissues are highly specific for single cell type Histogram of each mixture component's highest score for an annotated label. These highest scores as taken as the max of each column in Fig 3c,d and analogous cell annotation heatmaps for all tissues in SI Figure 3. Most mixture components score for a single annotated label uniquely.

SI Figure 5. Pairwise alignments between PopAlign models of Mammary Gland and Limb Muscle can be dissected in terms of Δw , $\Delta\mu$, $\Delta\Sigma$ Aligned subpopulations between the reference tissue, Mammary Gland, and the test tissue, Limb Muscle, are ranked by their Jeffrey's Divergence (a). For each aligned pair, we display: (b) associated p-value, (c) mean gene expression states (μ_i) in terms of annotated features, (d) shifts in abundance (Δw), (e) shifts in mean gene expression state ($\Delta\mu$), (f) shifts in population spread ($\Delta\Sigma$).

SI Figure 6. Pairwise alignments between PopAlign models of immune cells and signal conditions (IFNG, GM-CSF) can be dissected in terms of Δw , $\Delta\mu$, $\Delta\Sigma$ Aligned subpopulations between the reference PBMC sample and each test sample, GM-CSF and IFNG, are ranked by their Jeffrey's Divergence (a). For each aligned pair, we display: (b) mean gene expression states (μ_i) in terms of annotated features, (c) shifts in abundance (Δw), (d) shifts in mean gene expression state ($\Delta\mu$), (e) shifts in population spread ($\Delta\Sigma$).

SI Figure 7. Classification of reference patient population (healthy1) using canonical immune system markers. We classify cells using heatmap of mean gene expression values for cells from each reference subpopulation from healthy1. Each gene is independently scaled by its max value across all subpopulation means. Erythrocytes are HBB+, T cells are CD3D+, naive T cells are CCR+/CD62L+/CD27+/CD28+, effector T cells are CD57+/PD1+/CD95+, B cells are CD19+, and monocytes are CD14+/CD16+.

SI Figure 8. Pairwise alignments between PopAlign models of healthy1 and healthy2 Aligned subpopulations between the reference healthy1 sample and test sample, healthy2, are ranked by their Jeffrey's Divergence. For each aligned pair, the following are displayed: mean gene expression states (μ_i) in terms of features, shifts in abundance (Δw), shifts in mean gene expression state ($\Delta\mu$), shifts in population spread ($\Delta\Sigma$), and $\log_{10}(\text{Jeffrey's})$

Divergence).

SI Figure 9. Pairwise alignments between PopAlign models of healthy1 and MM1 See caption for SI Figure 8.

SI Figure 10. Pairwise alignments between PopAlign models of healthy1 and MM2 See caption for SI Figure 8.

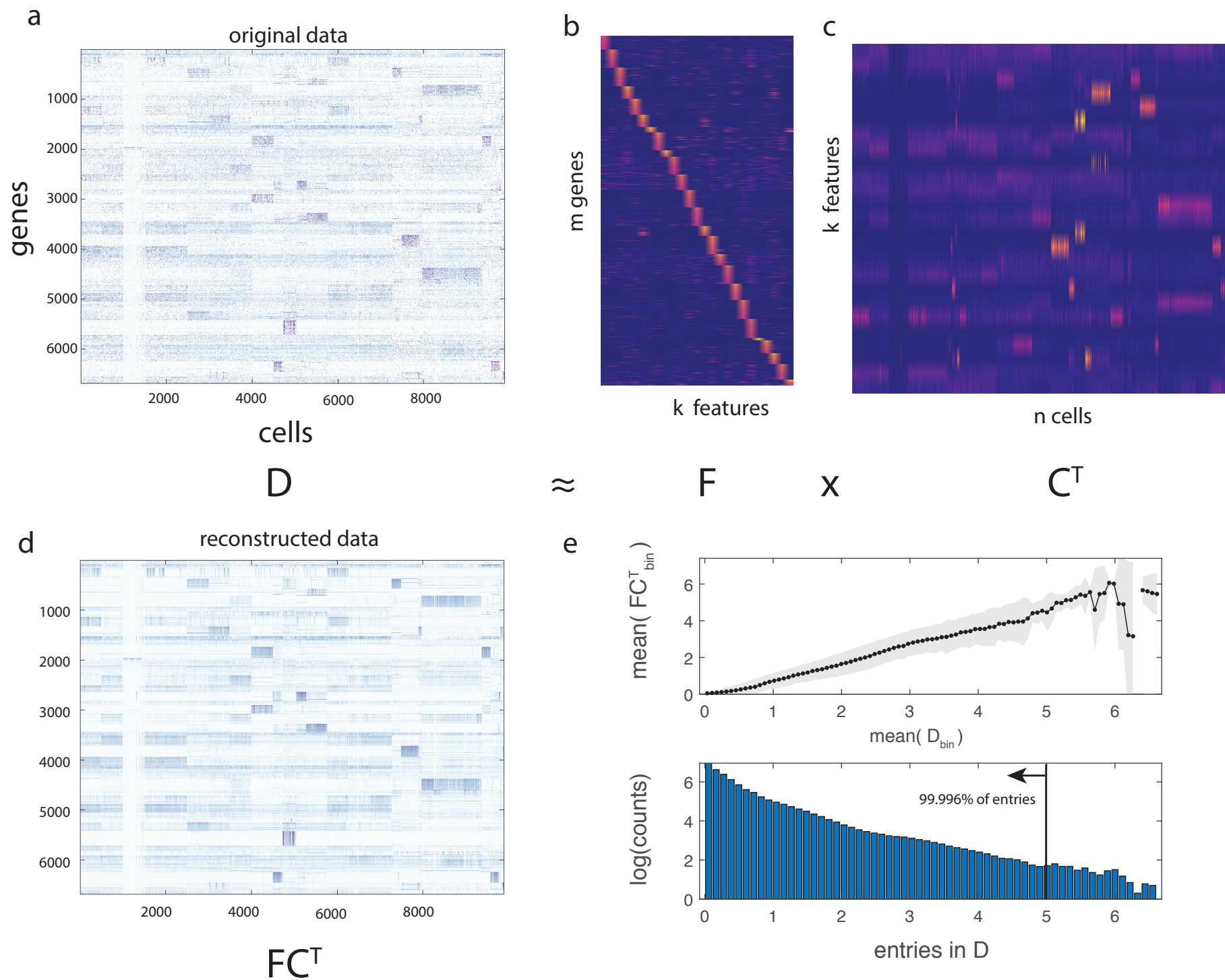
SI Figure 11. Pairwise alignments between PopAlign models of healthy1 and MM3 See caption for SI Figure 8.

SI Figure 12. Pairwise alignments between PopAlign models of healthy1 and MM4 See caption for SI Figure 8.

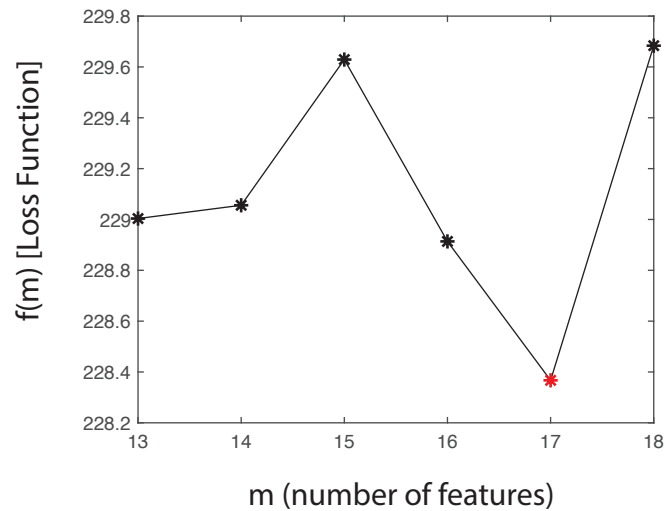
SI Figure 13. t-SNE analysis obscures the relationships between aligned subpopulations (a-c) tSNE plots of single cell data for three signaling conditions from Figure 4: (a) PBMC, (b) GM-CSF, and (c) IFNG. Mixture component μ 's (disks) for PBMC mixture components 1 and 3 overlaid on plots (b) and (c) to aid interpretation. PopAlign generated alignments for monocyte mixtures from figure (3) are indicated as lines. tSNE was performed on data pooled from all experiments, but experimental conditions are shown individually for clarity. Plot shows that gene expression changes due to signaling can result in relative cluster shifts within the tSNE plots that obscure subpopulation alignment. For example, in (b) the monocyte populations aligned by PopAlign (Figure 4) become separated by a cluster of T-cells (mixture components 3 and 4).

SI Table 1. Patient information for PBMC donors Patient information table with Age, Gender, Ethnicity, Disease Status, and Medications.

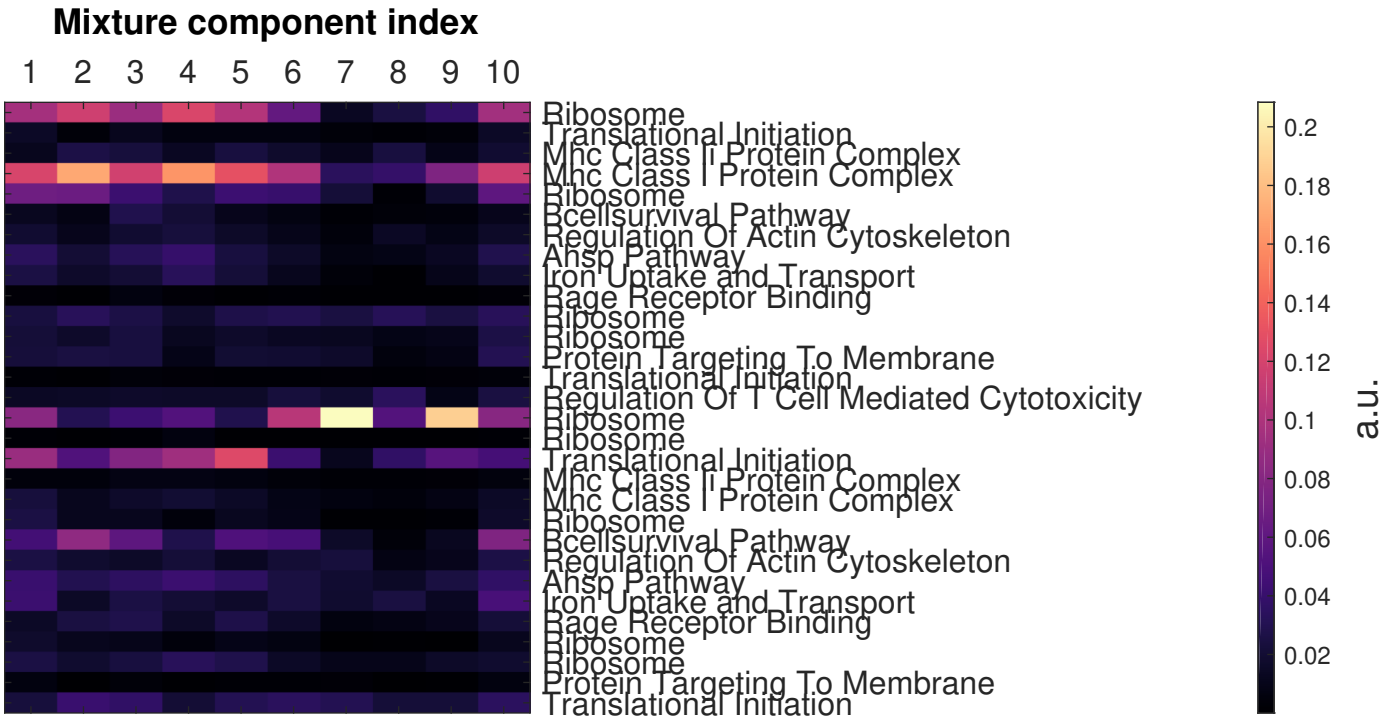
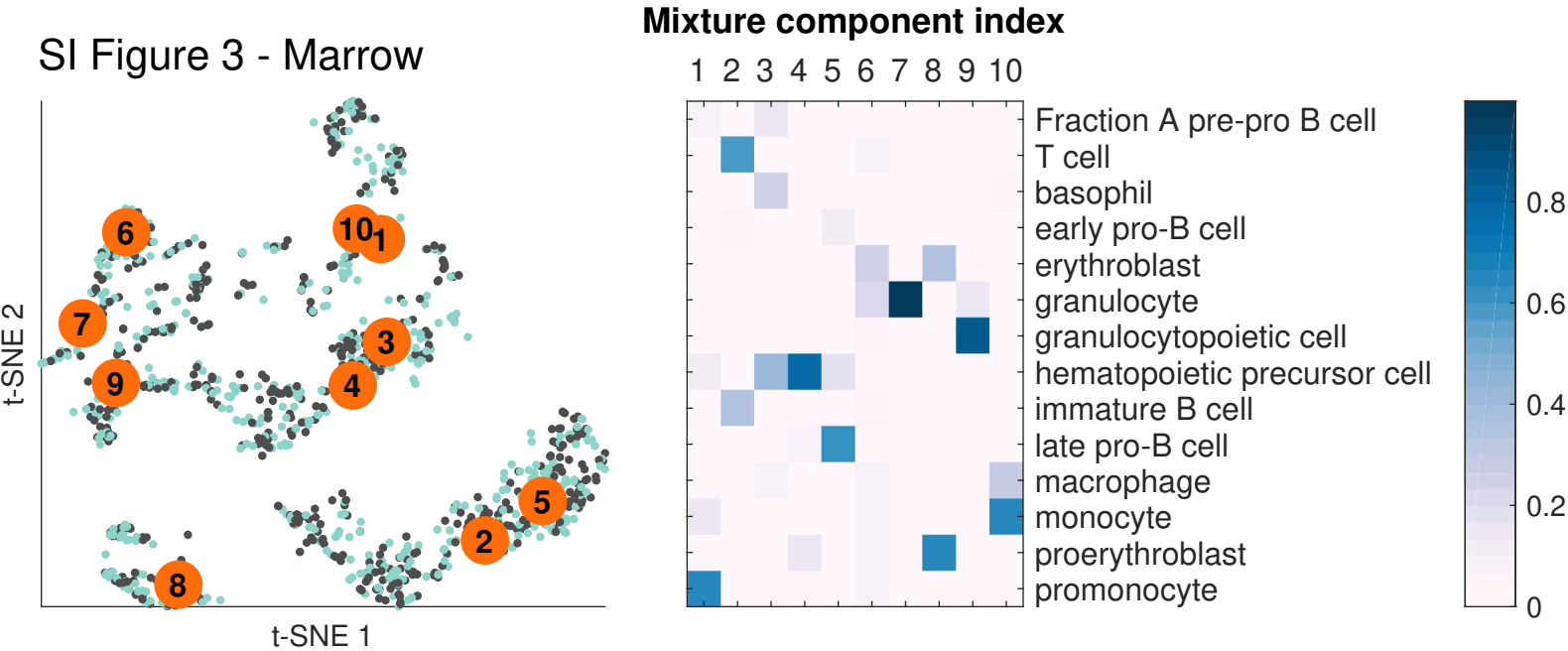
SI Figure 1



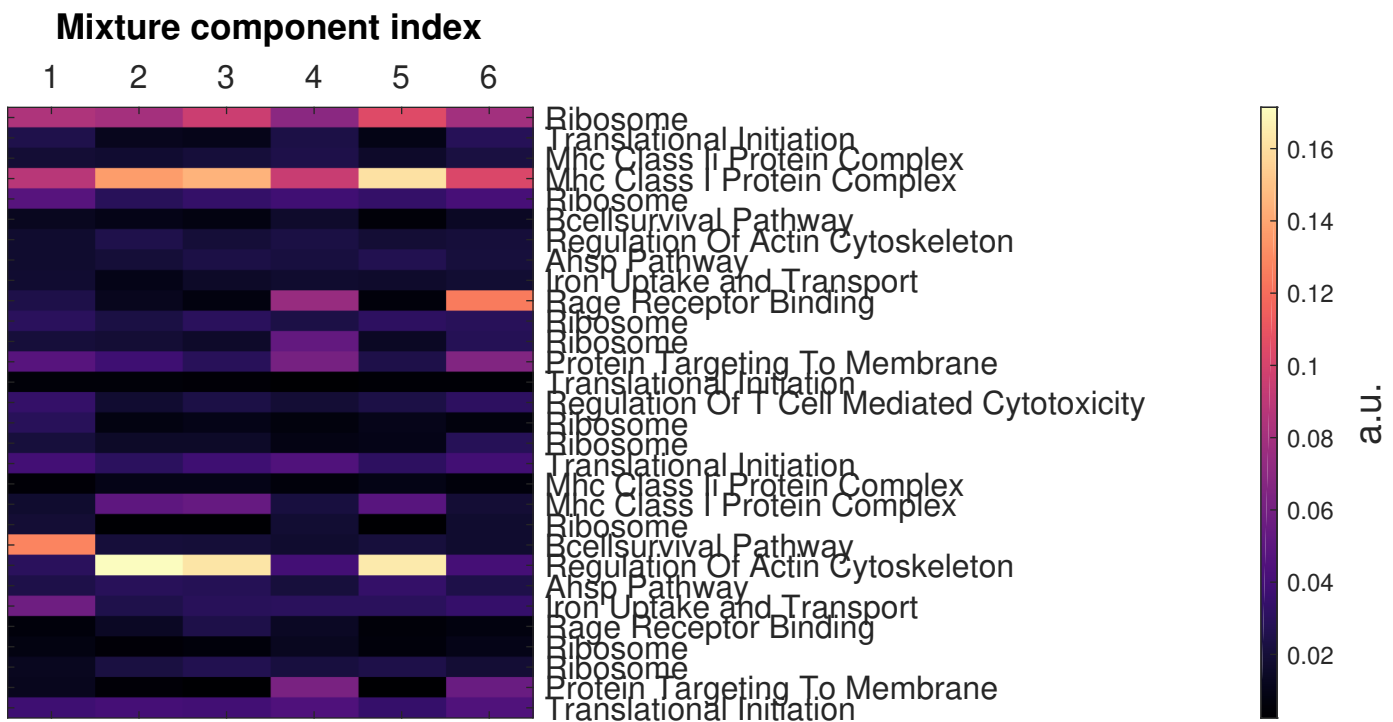
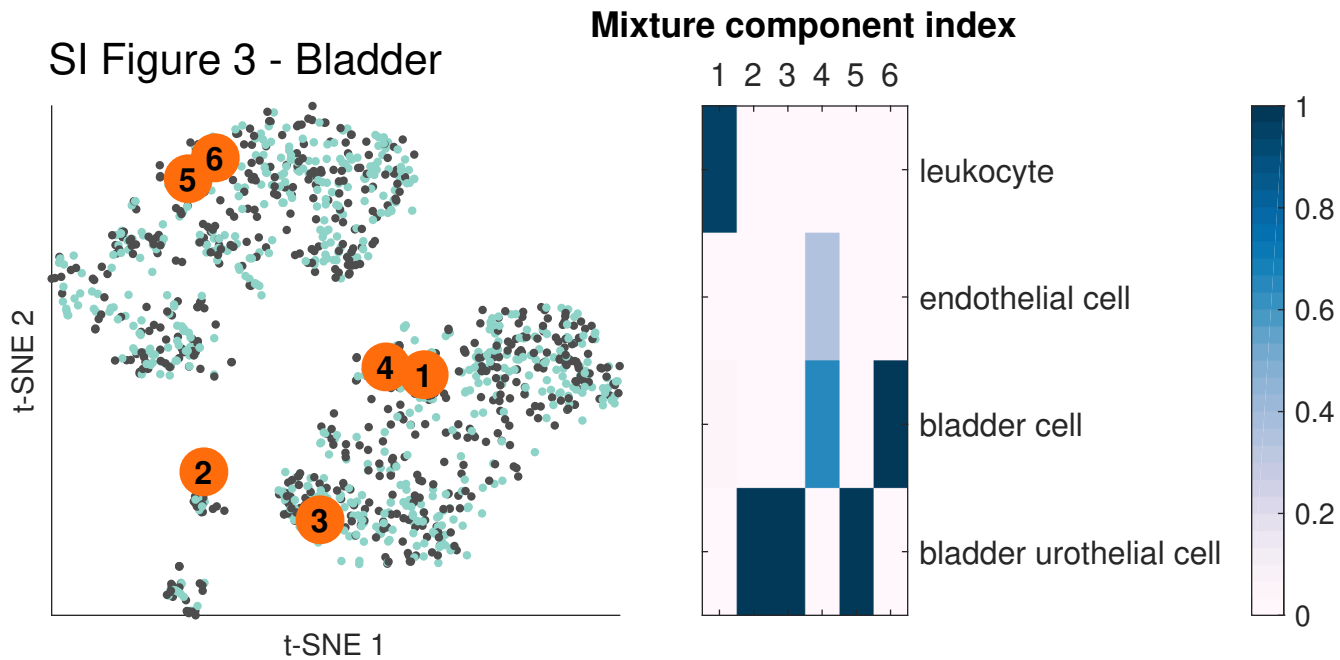
SI Figure 2



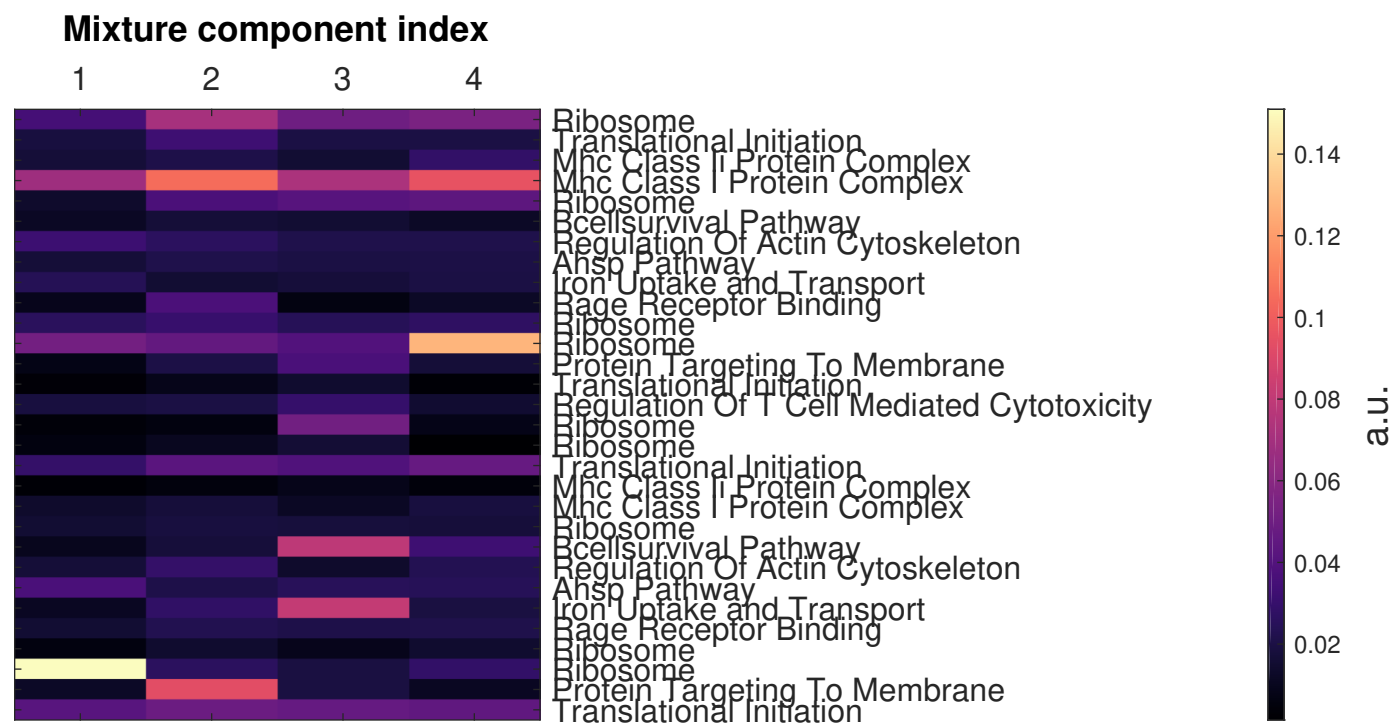
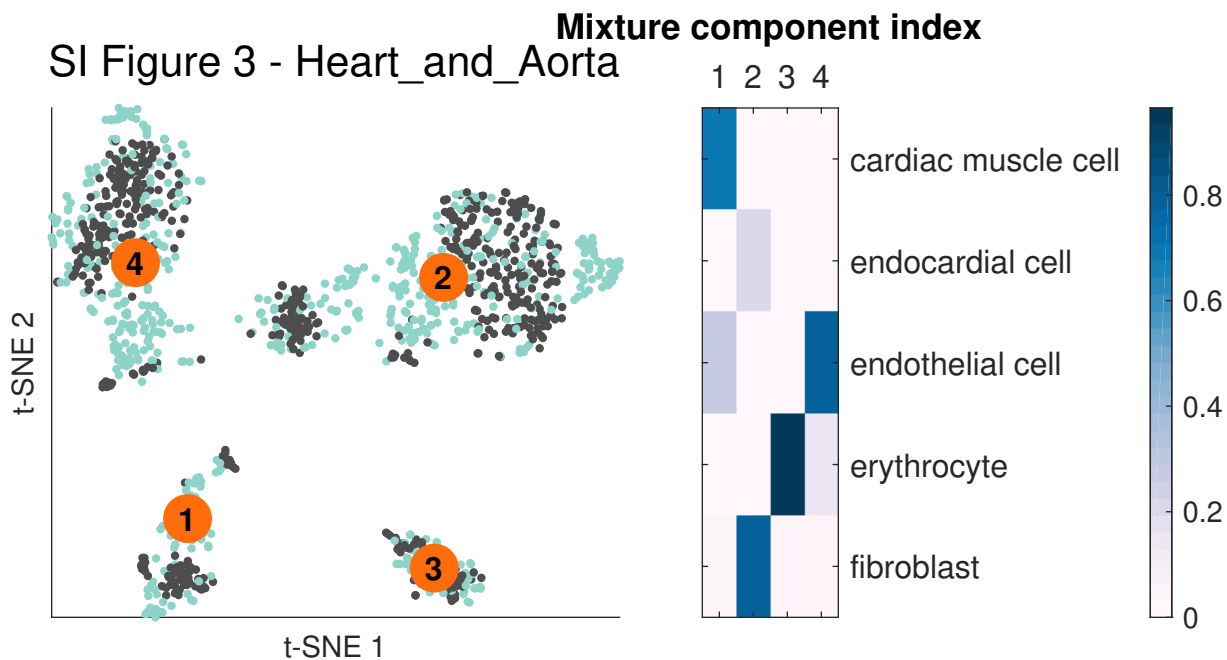
SI Figure 3 - Marrow



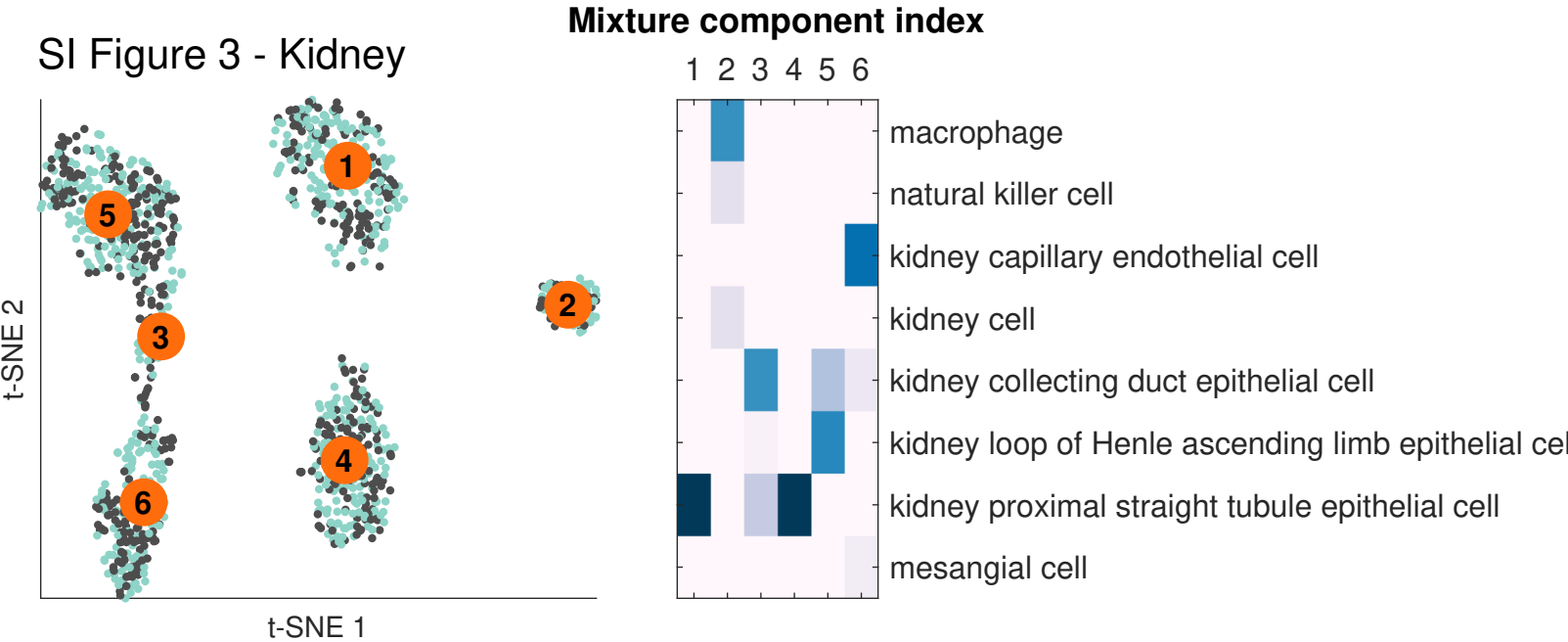
SI Figure 3 - Bladder



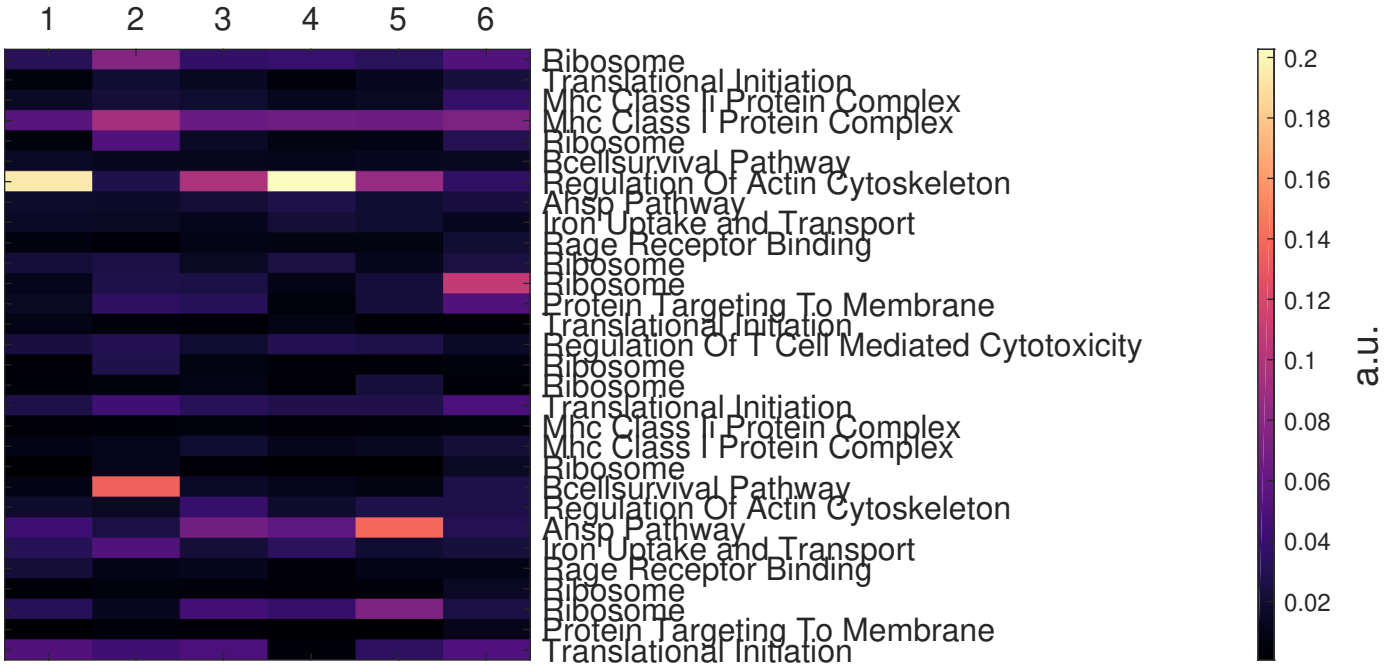
SI Figure 3 - Heart_and_Aorta



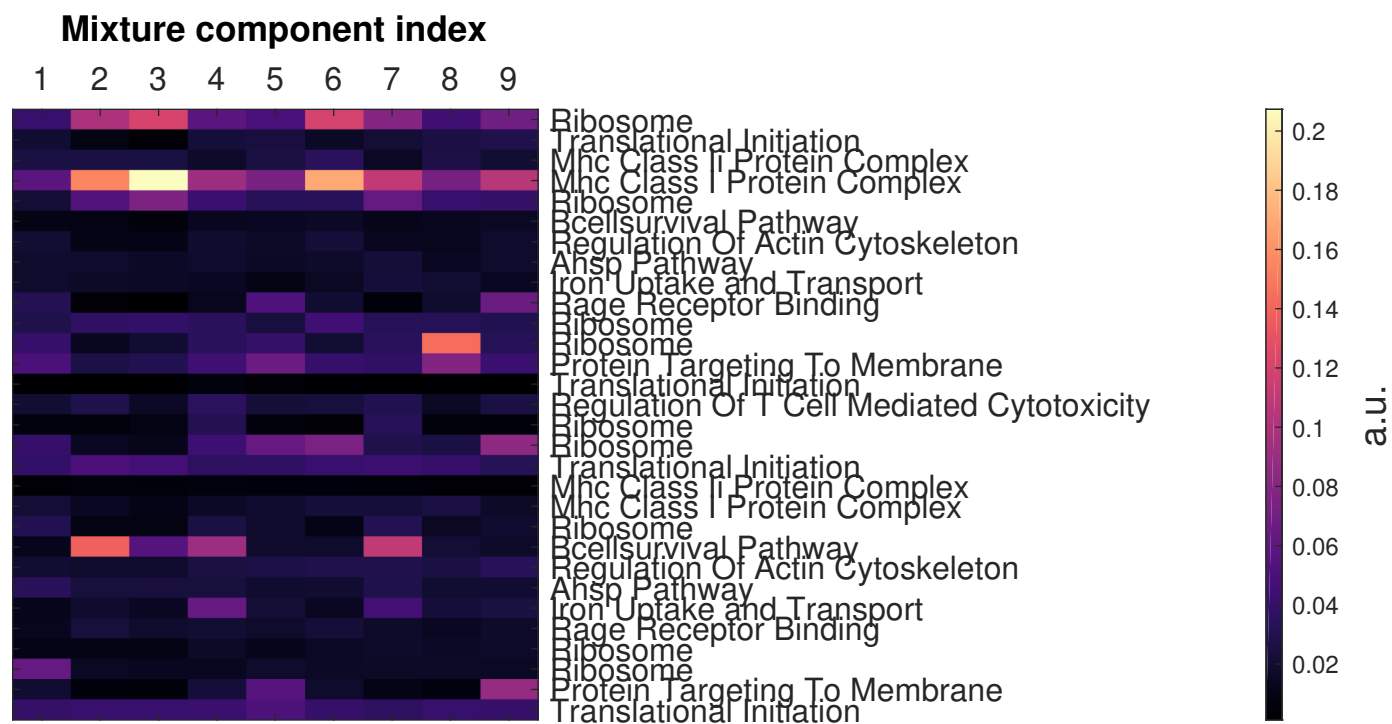
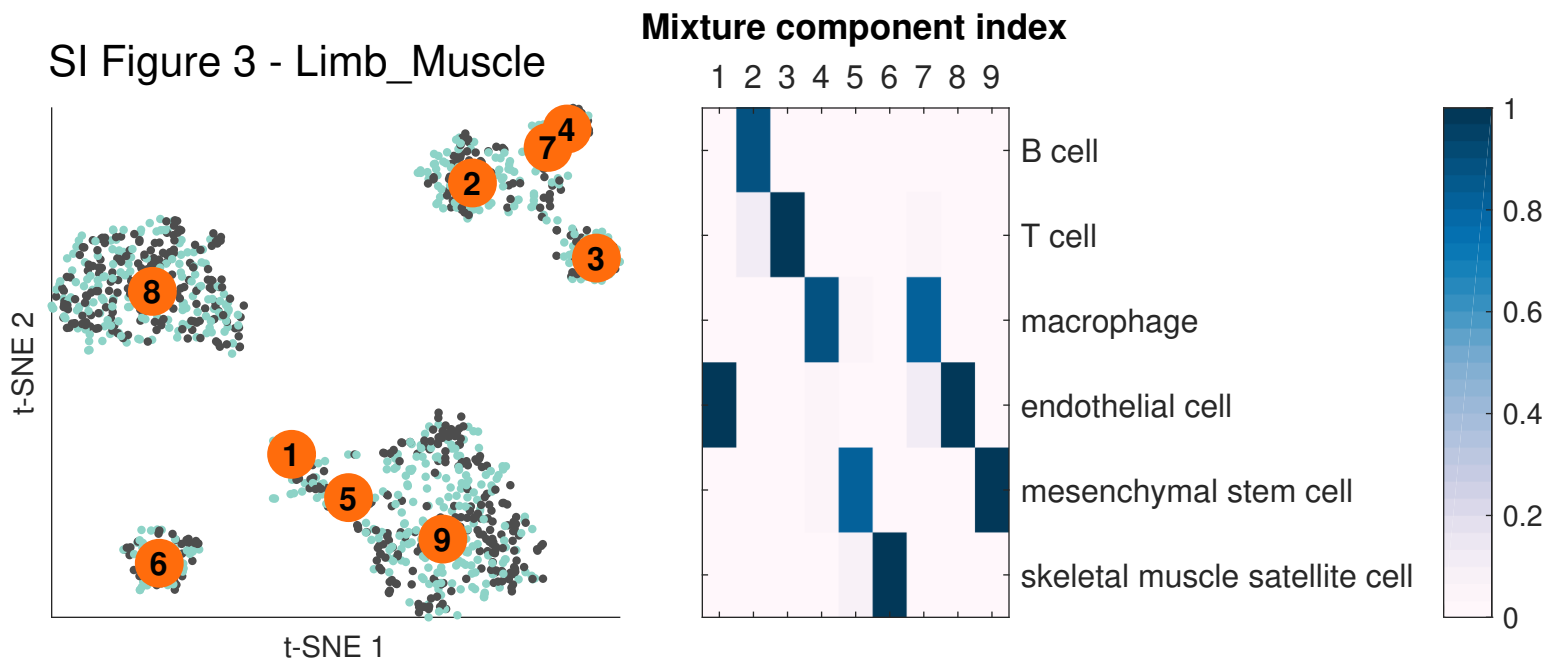
SI Figure 3 - Kidney



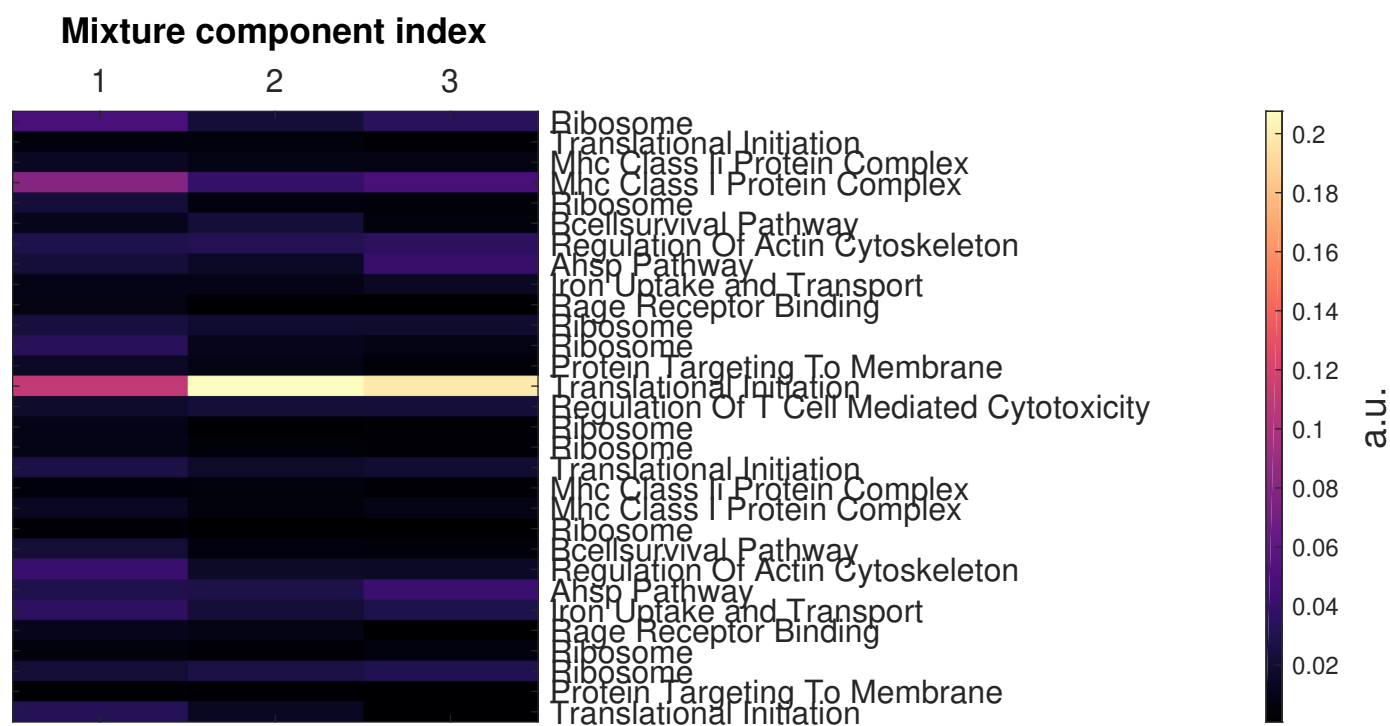
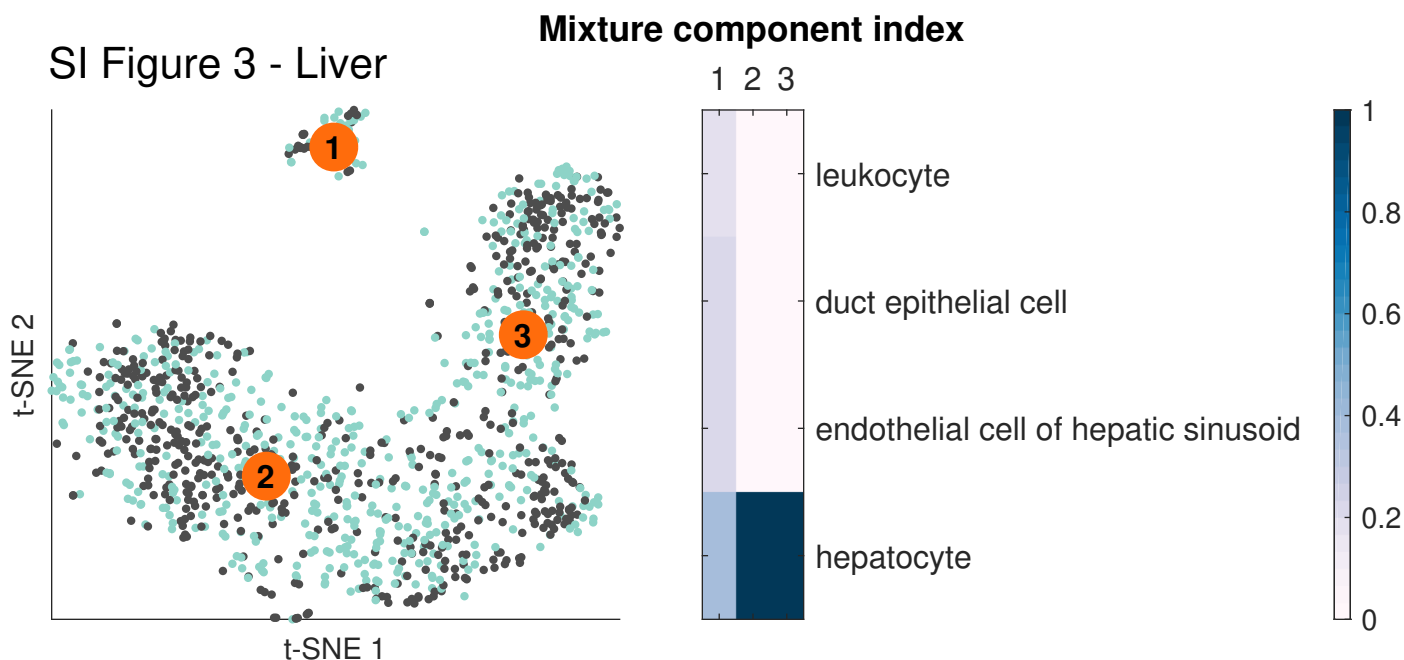
Mixture component index



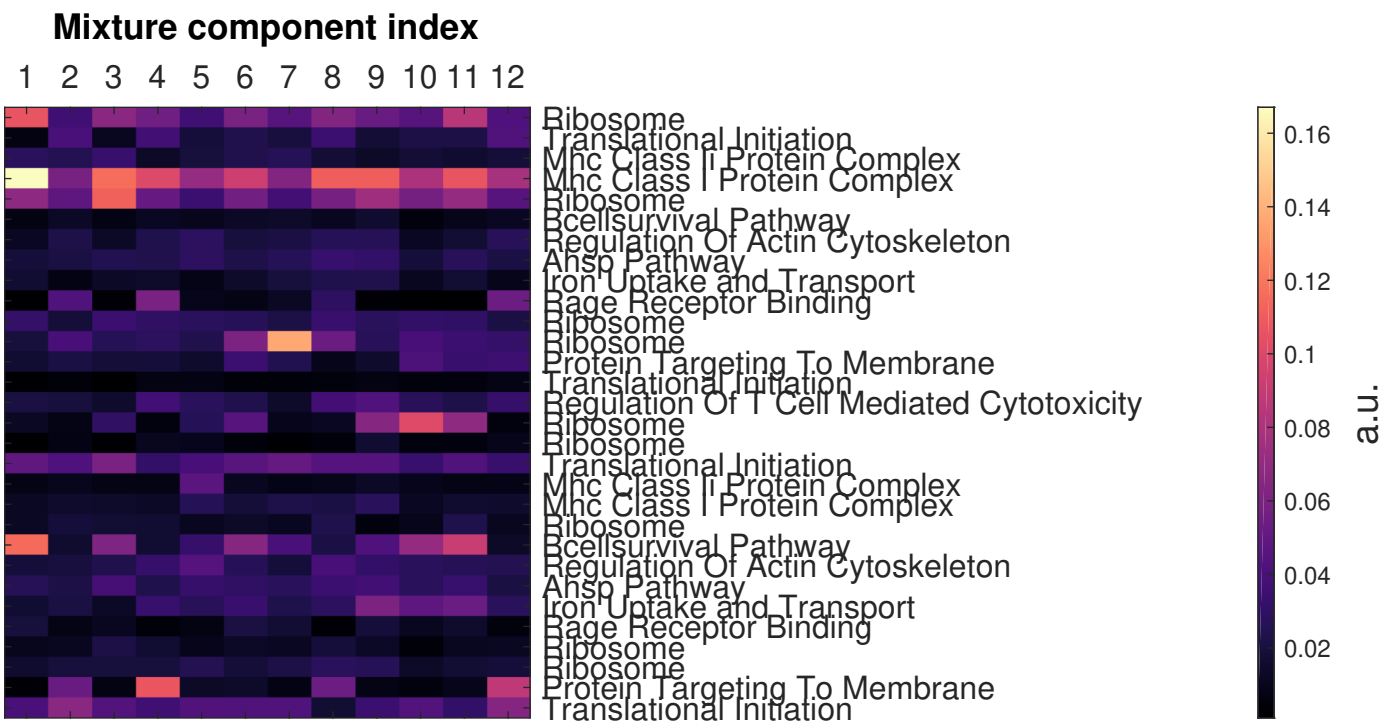
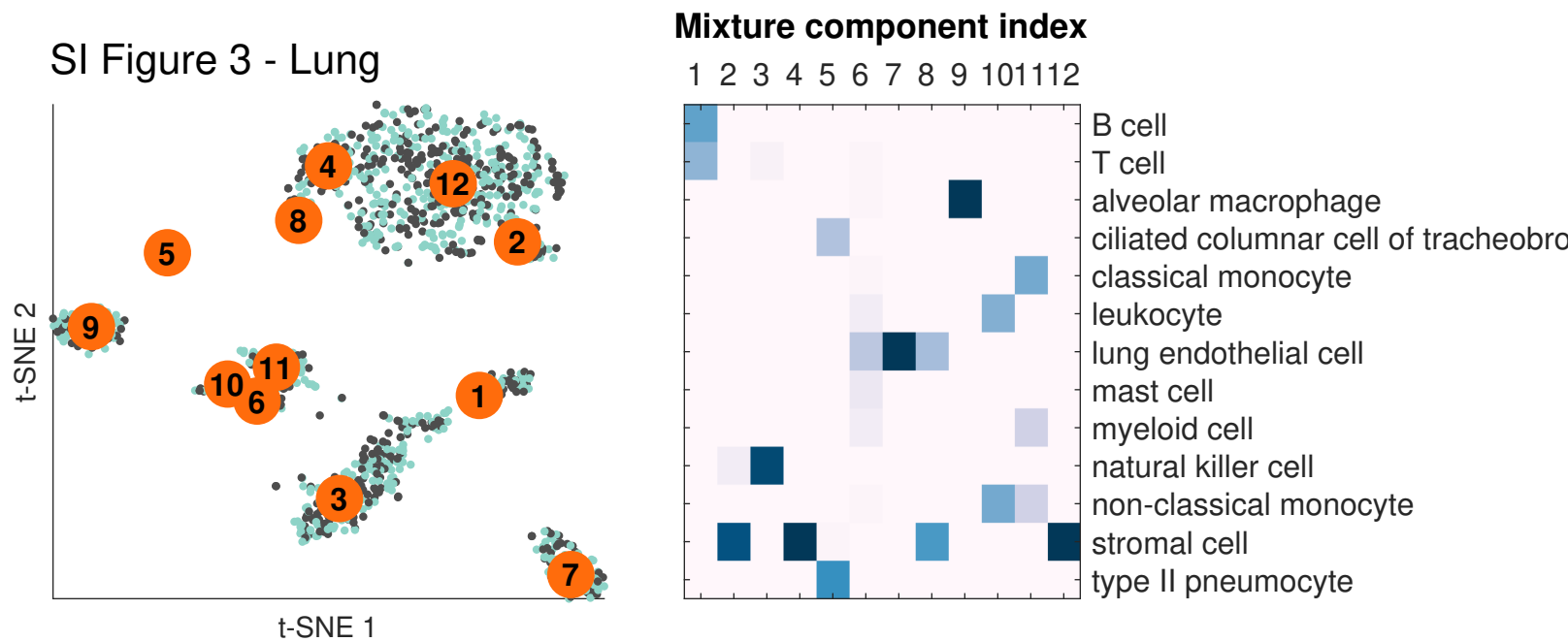
SI Figure 3 - Limb_Muscle



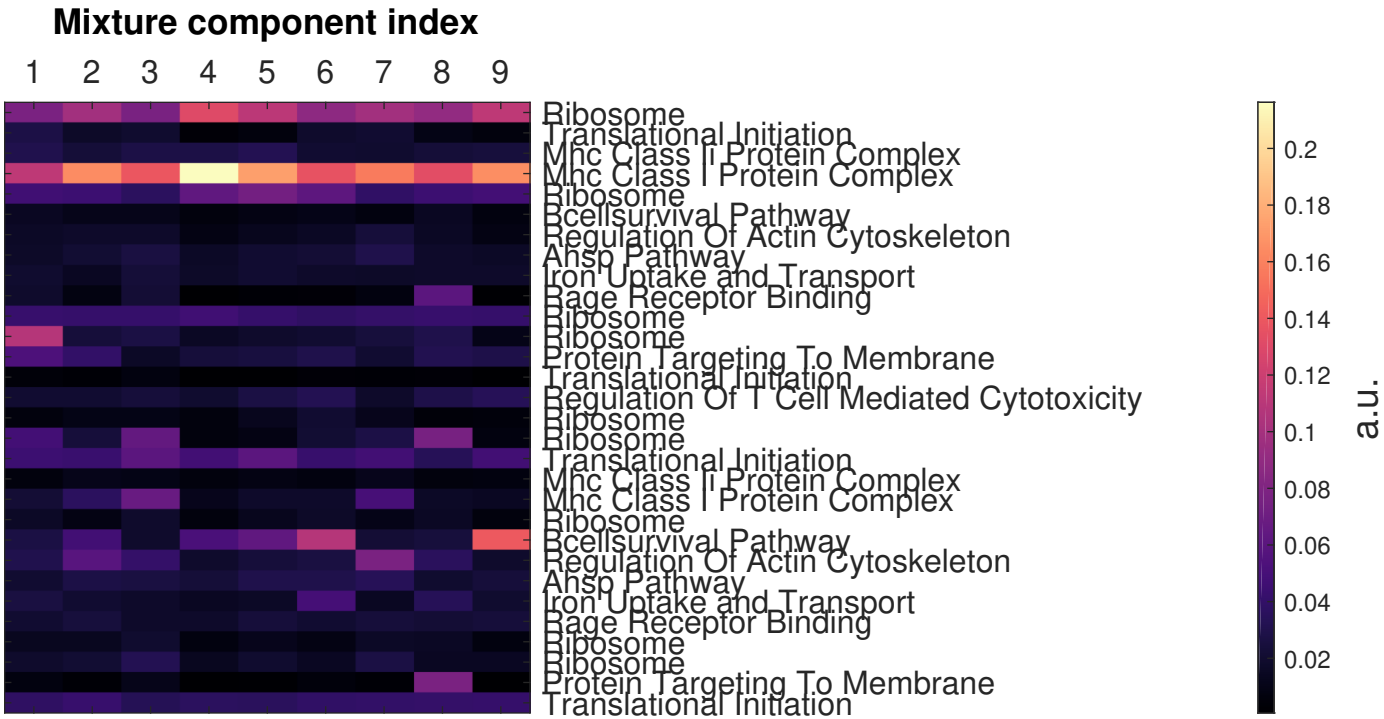
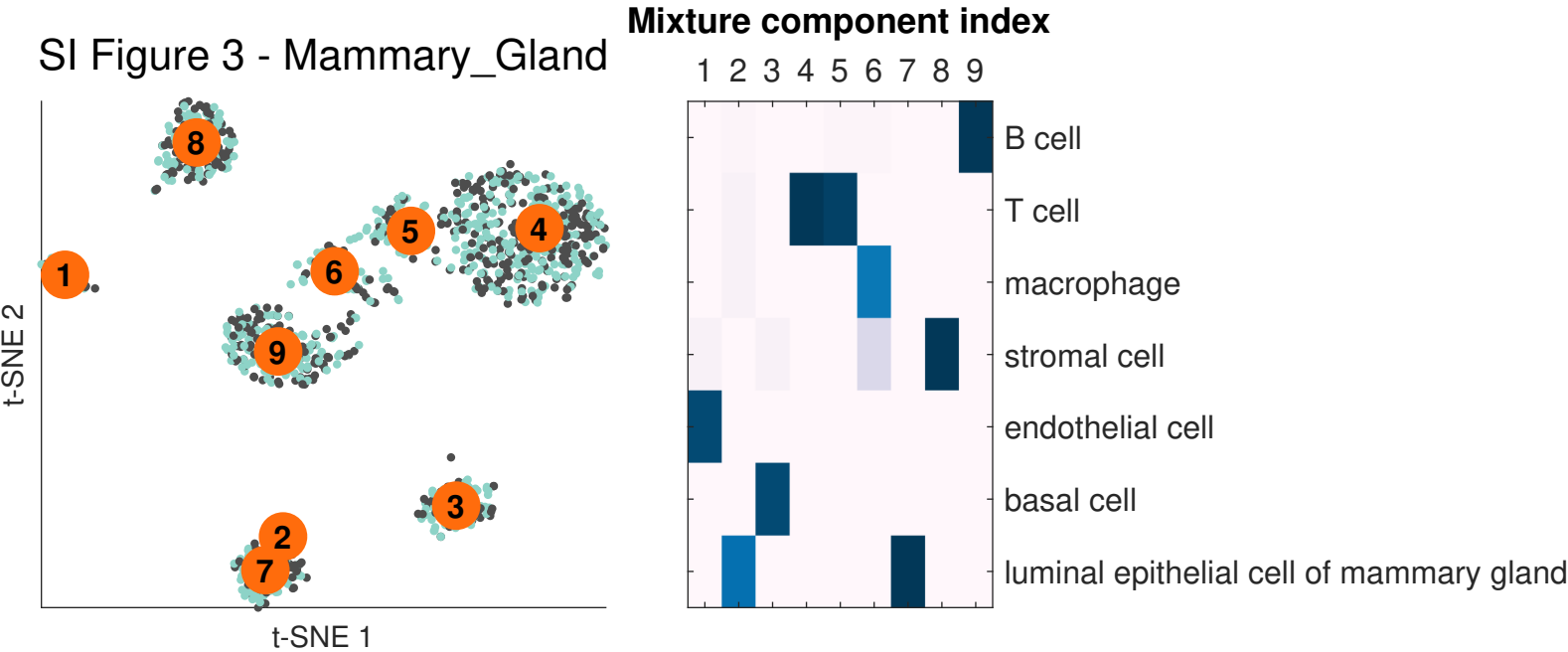
SI Figure 3 - Liver



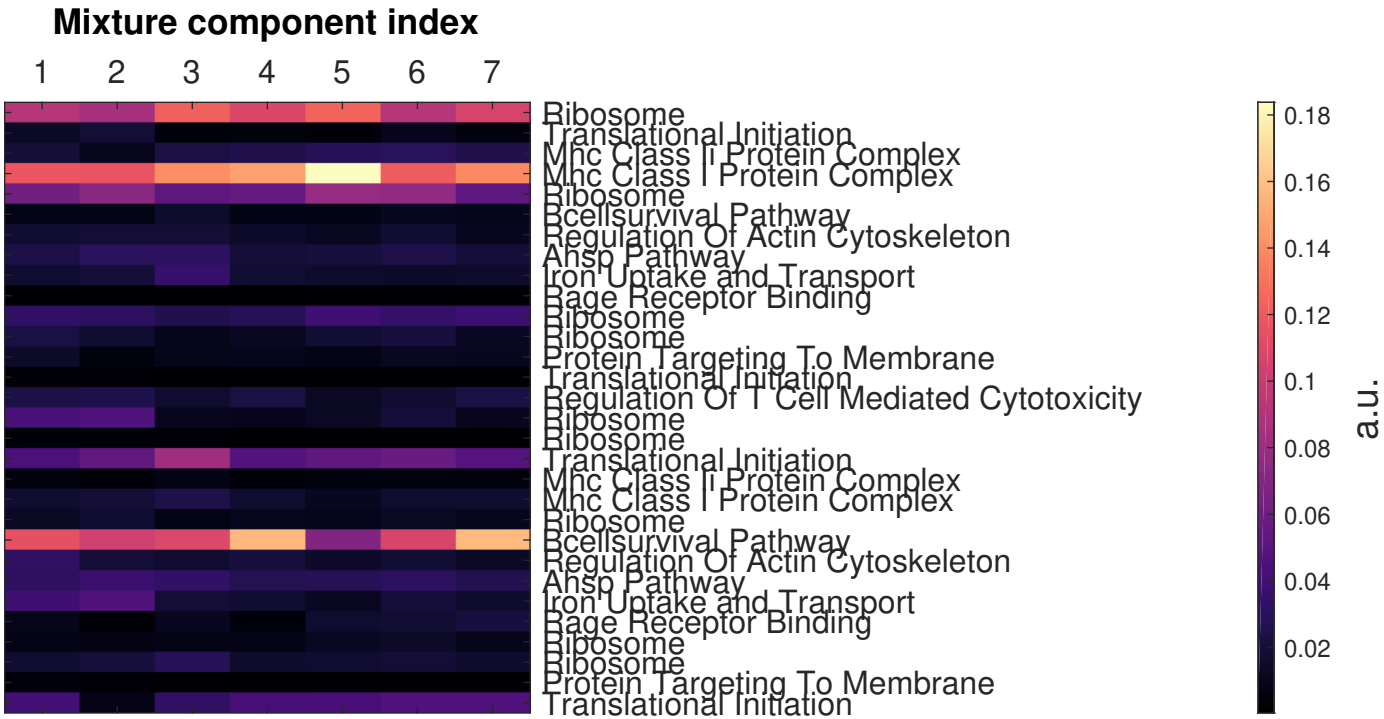
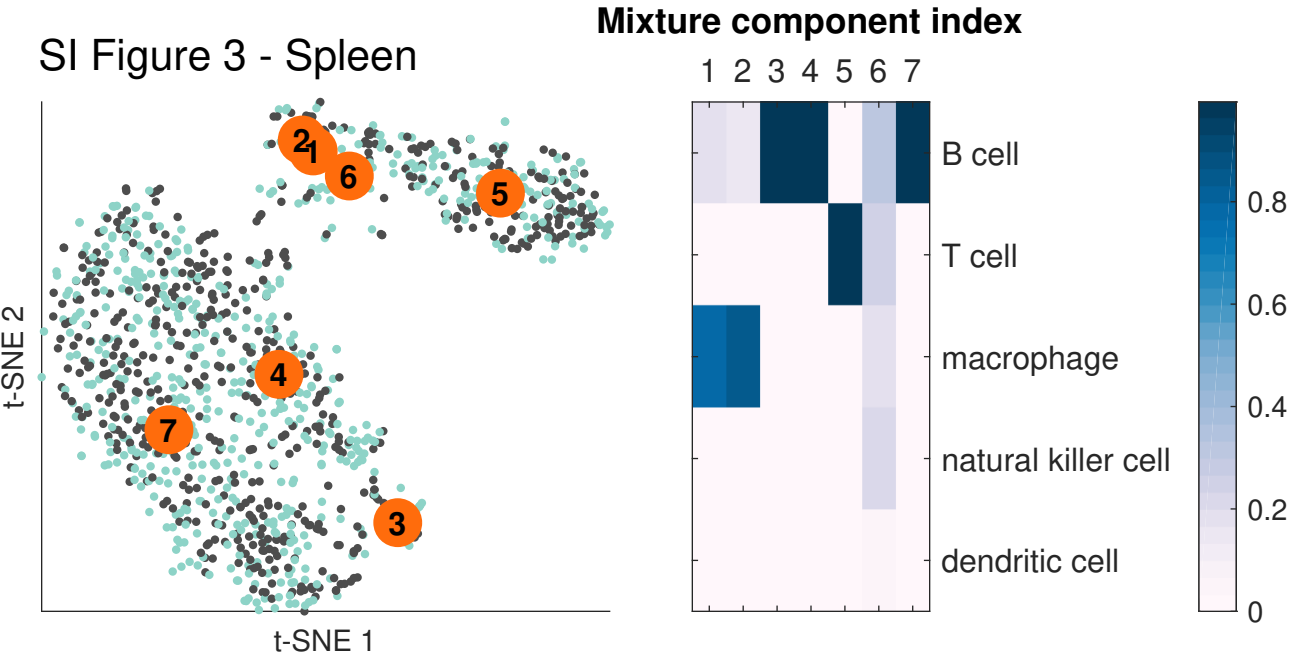
SI Figure 3 - Lung



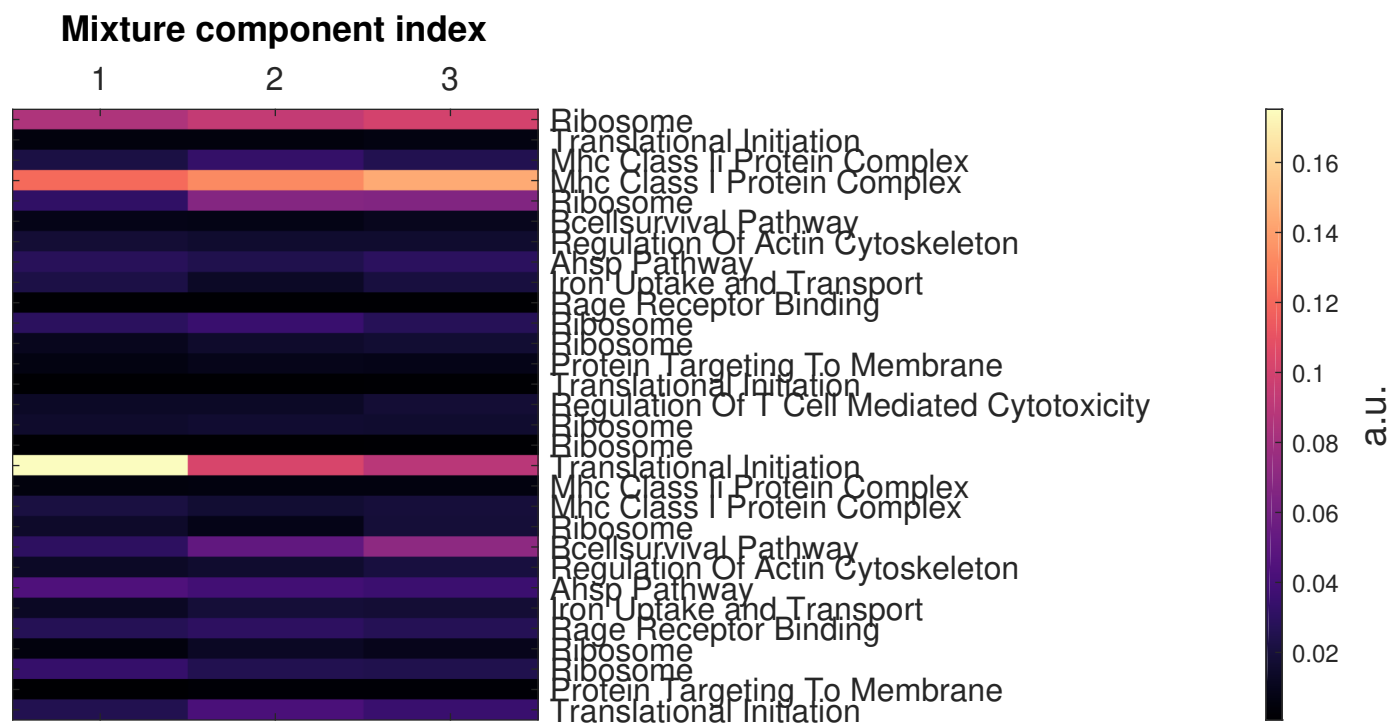
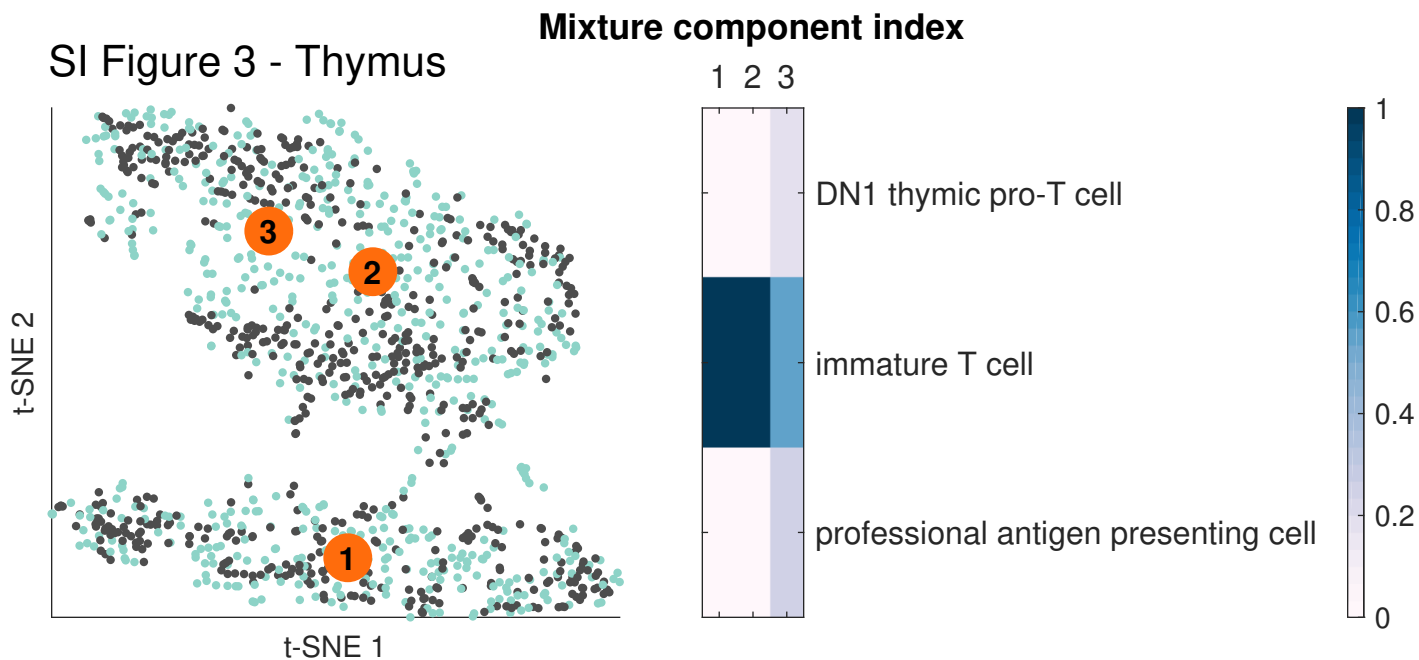
SI Figure 3 - Mammary_Gland



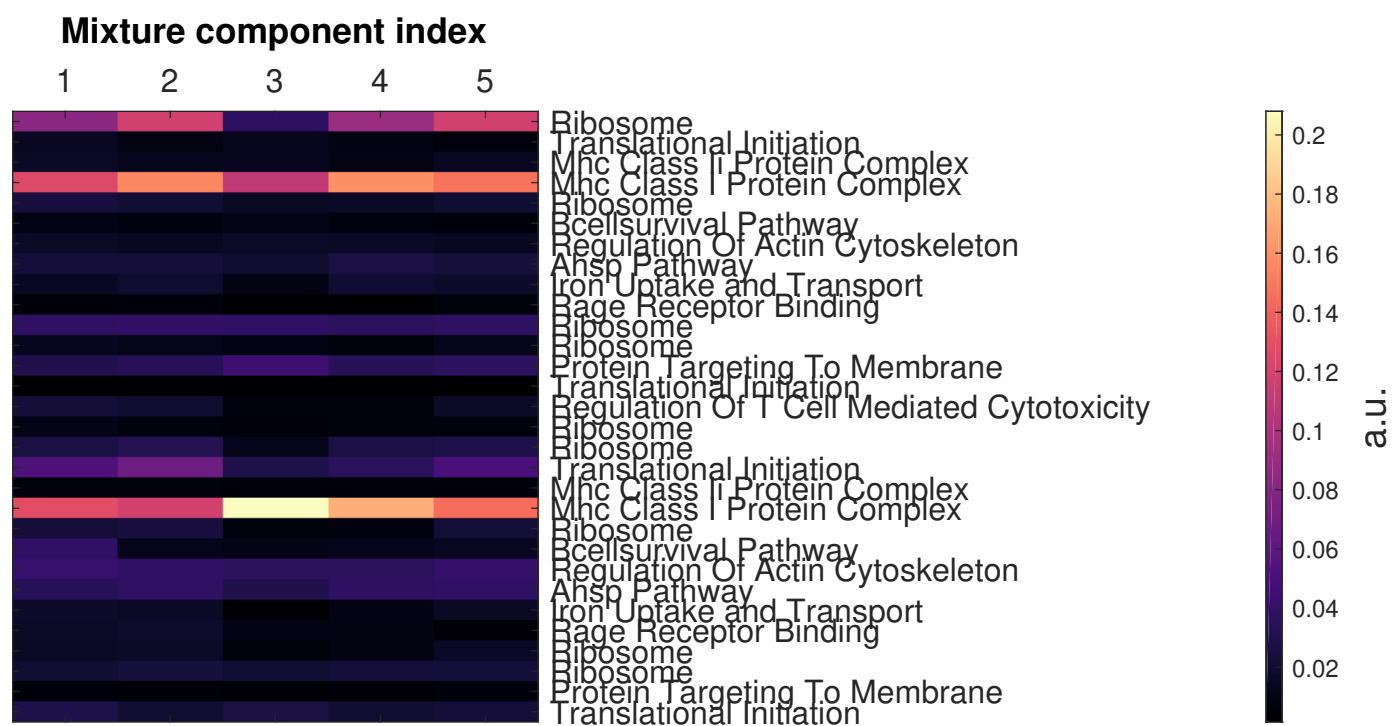
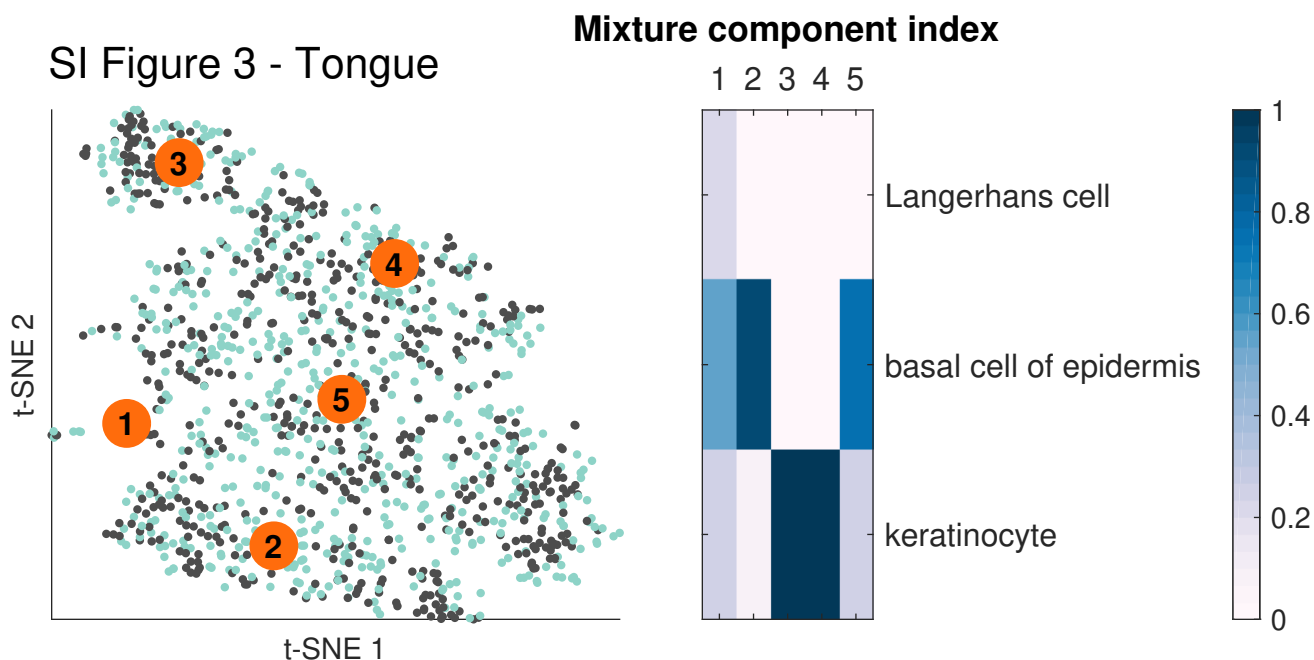
SI Figure 3 - Spleen



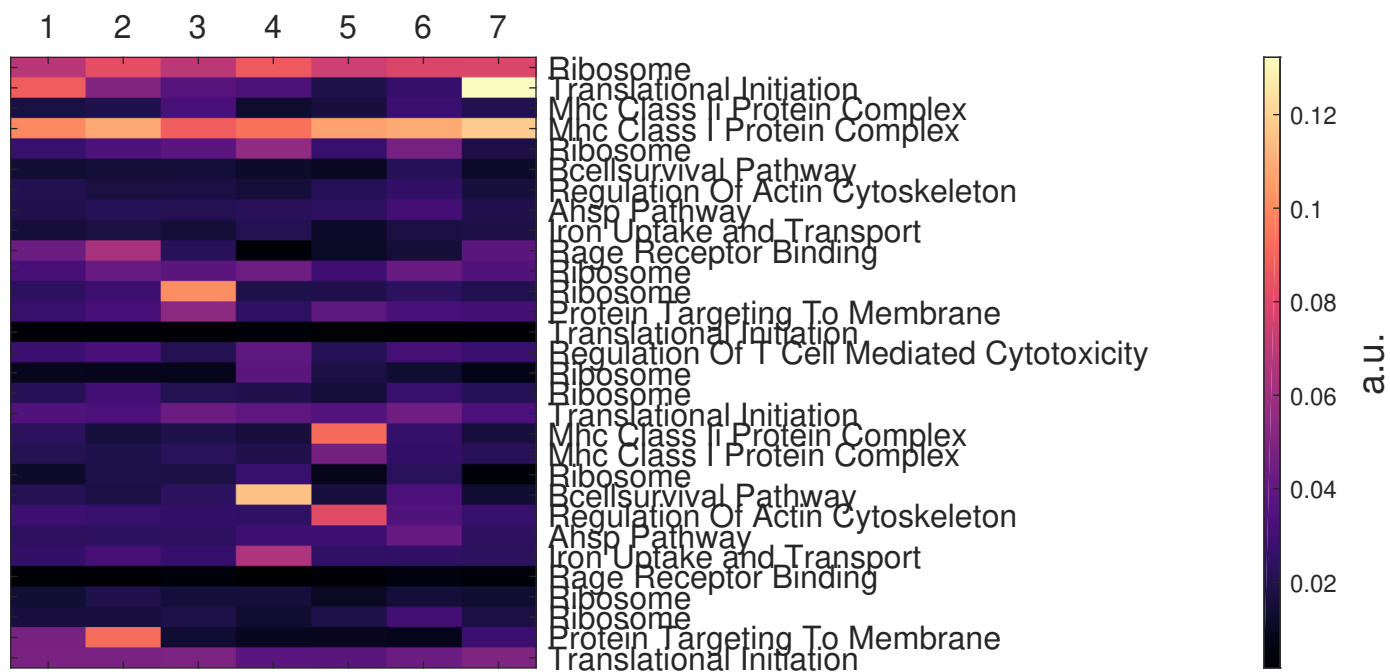
SI Figure 3 - Thymus



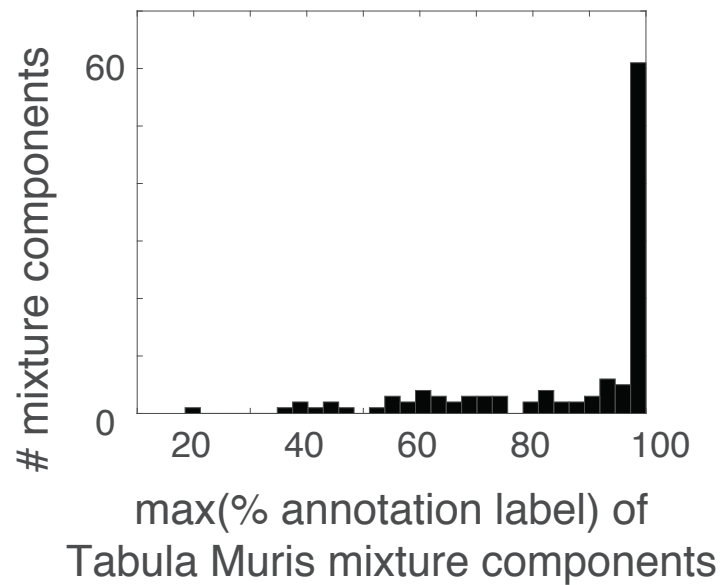
SI Figure 3 - Tongue



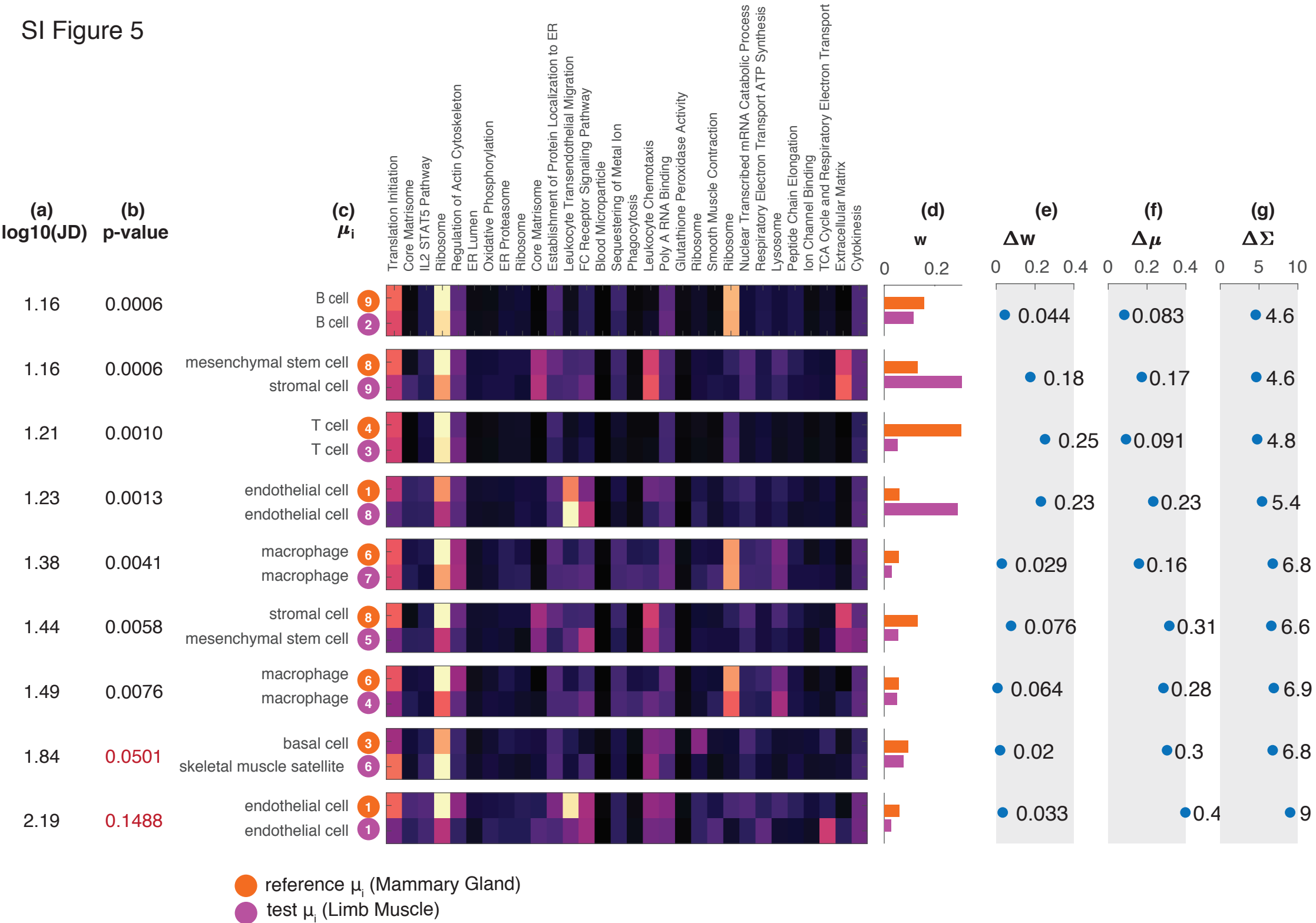
Mixture component index



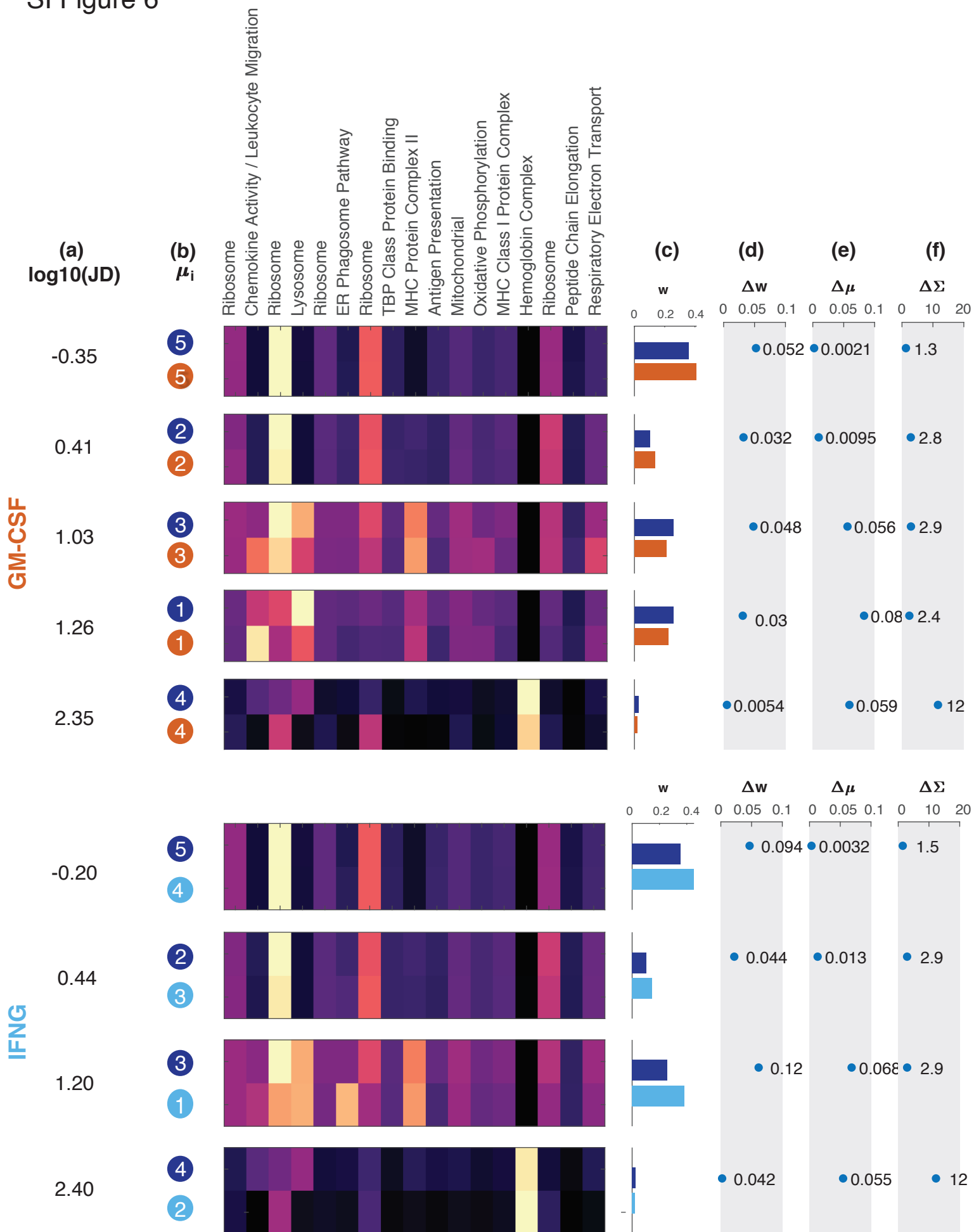
SI Figure 4



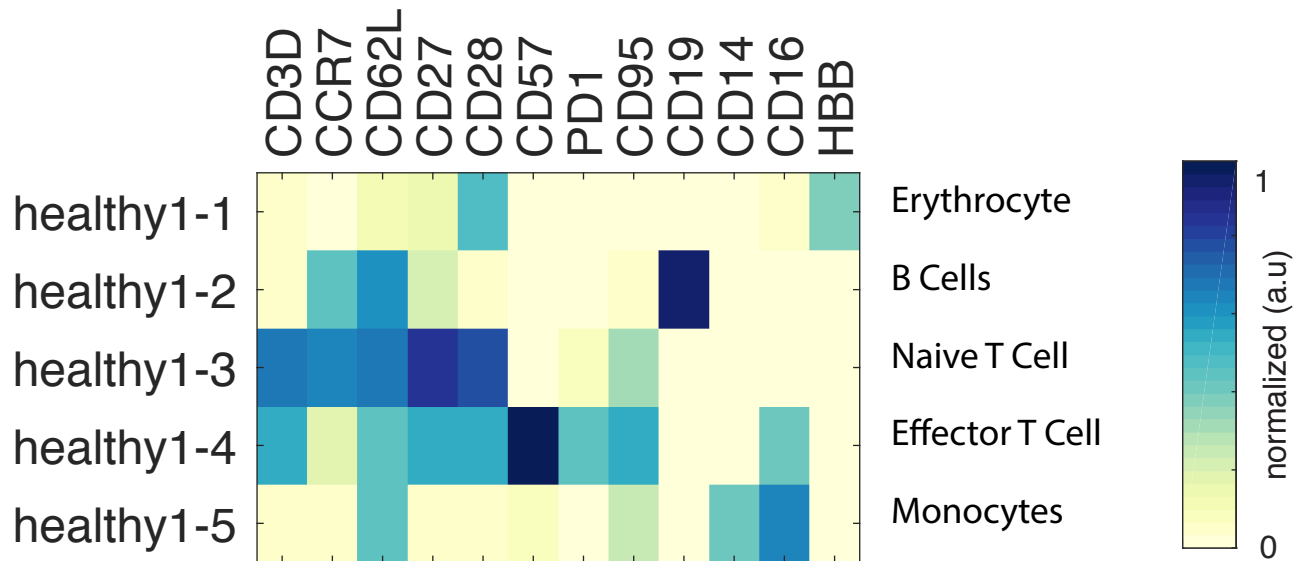
SI Figure 5



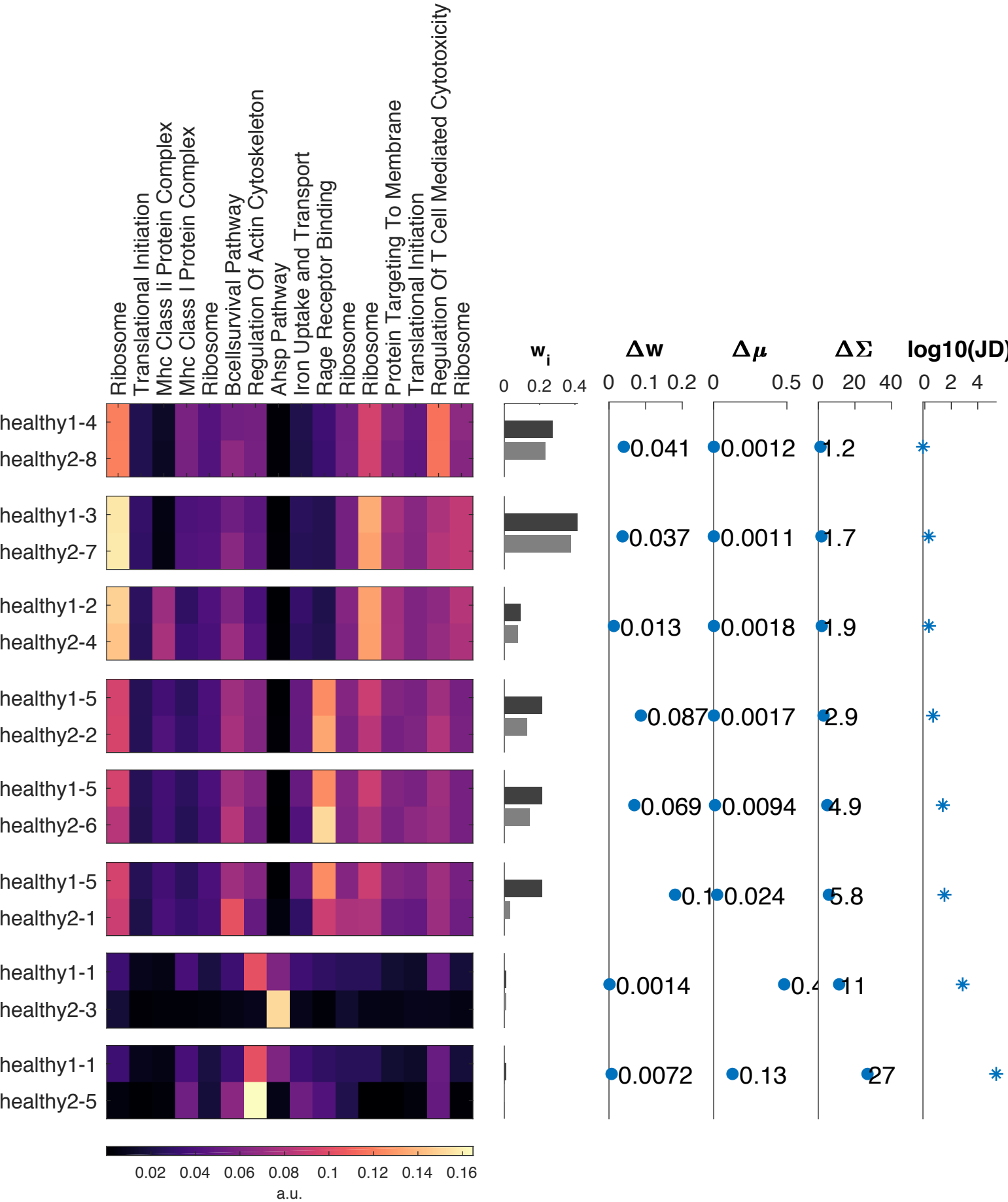
SI Figure 6



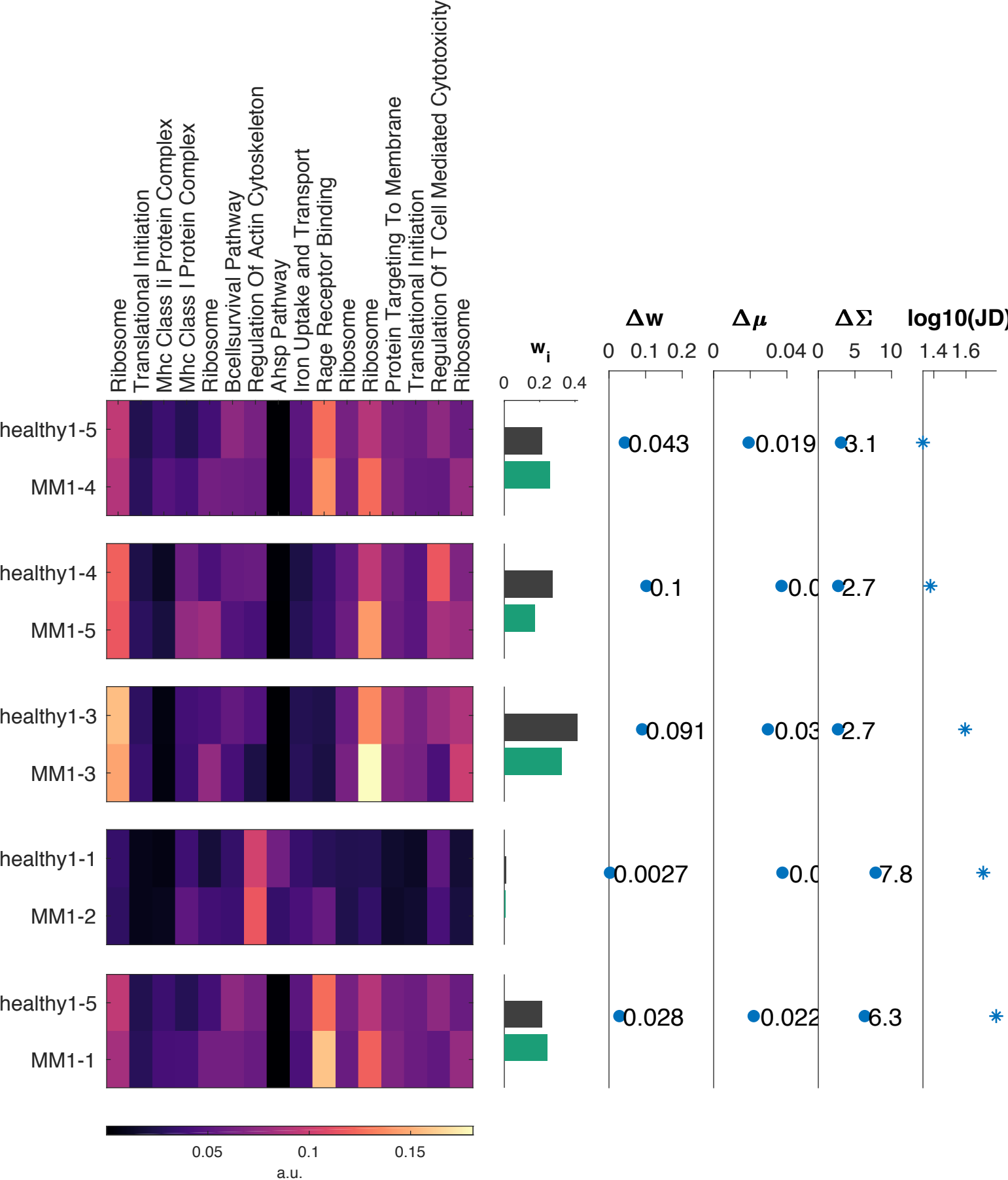
SI Figure 7



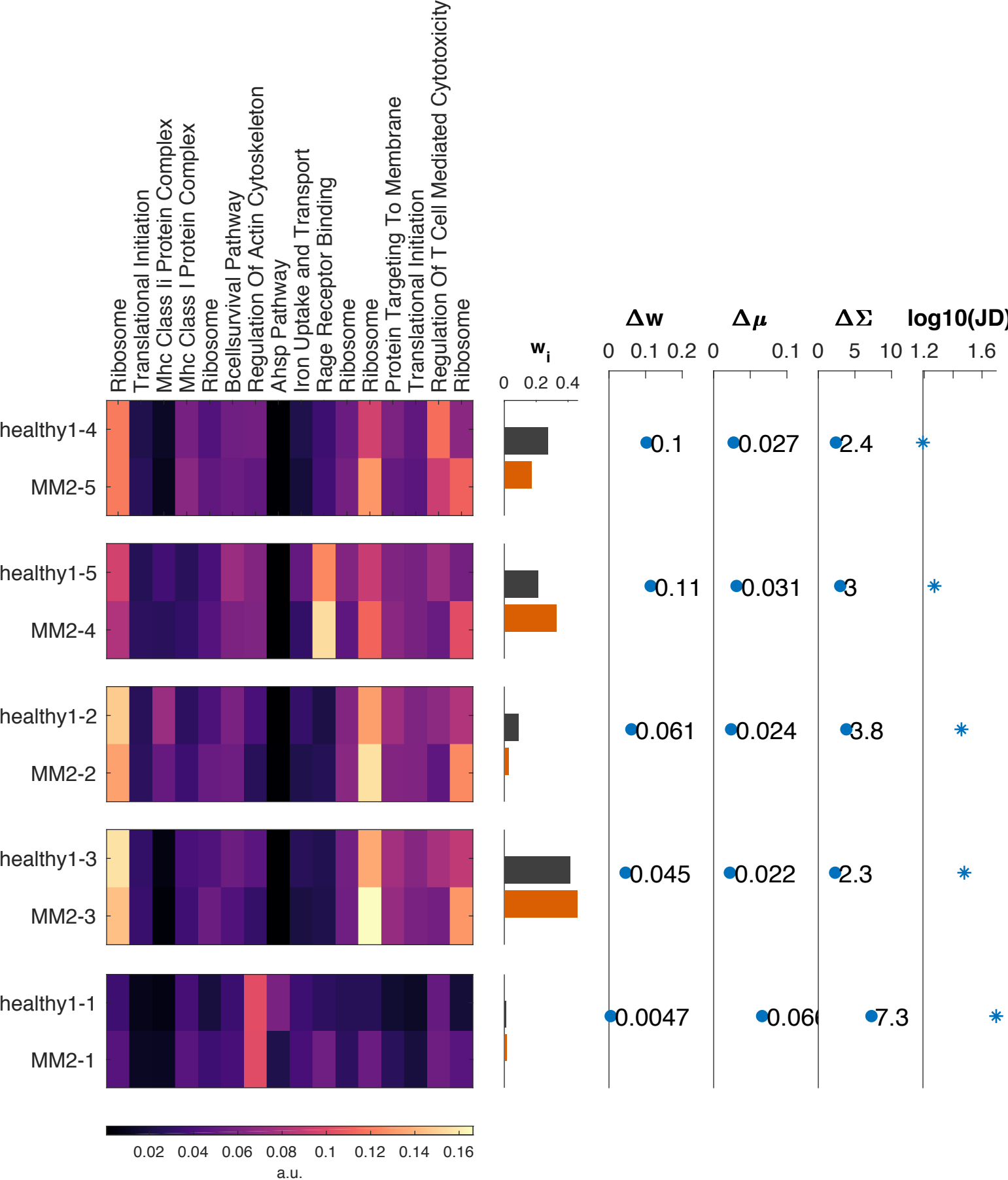
SI Figure 8



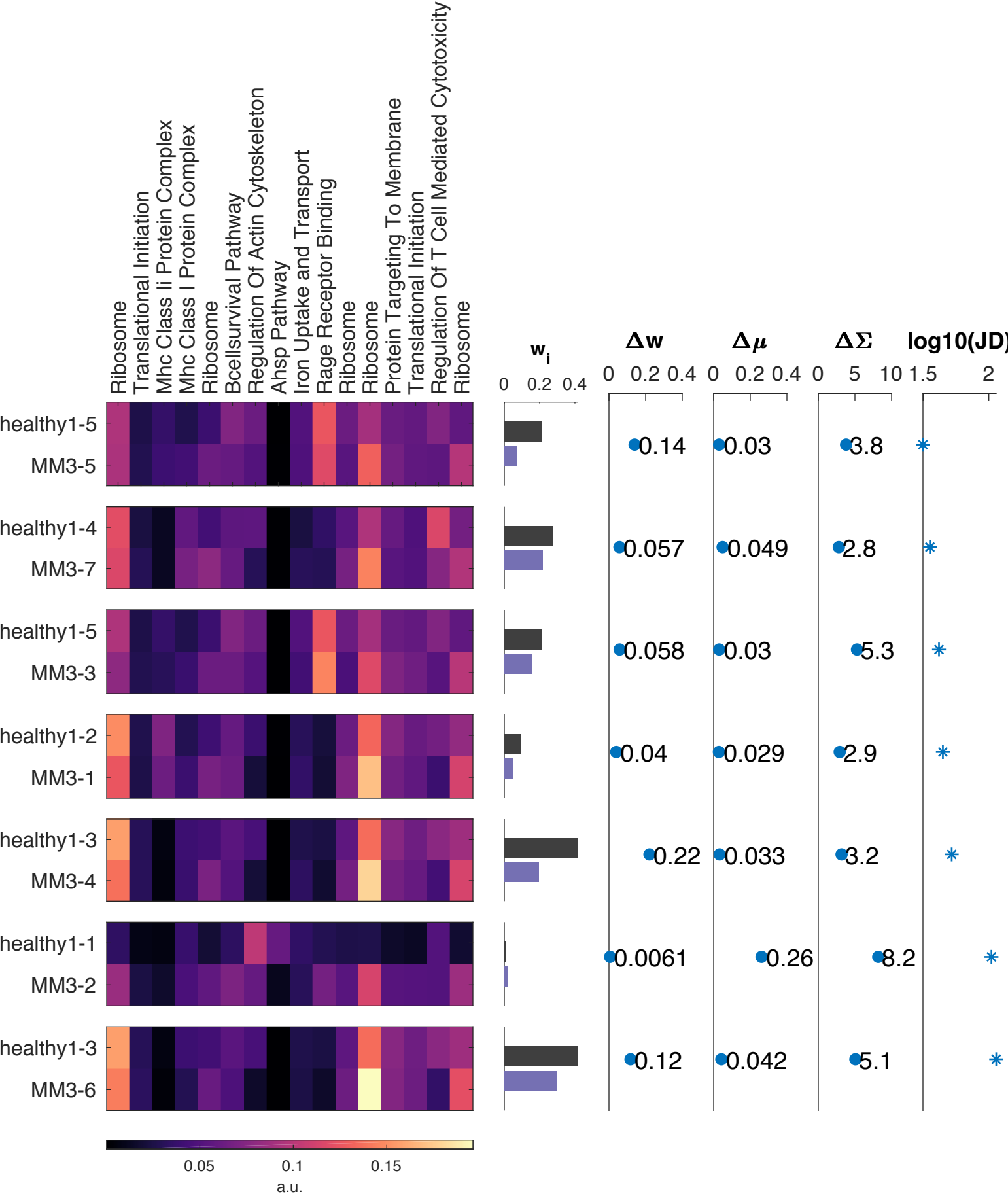
SI Figure 9



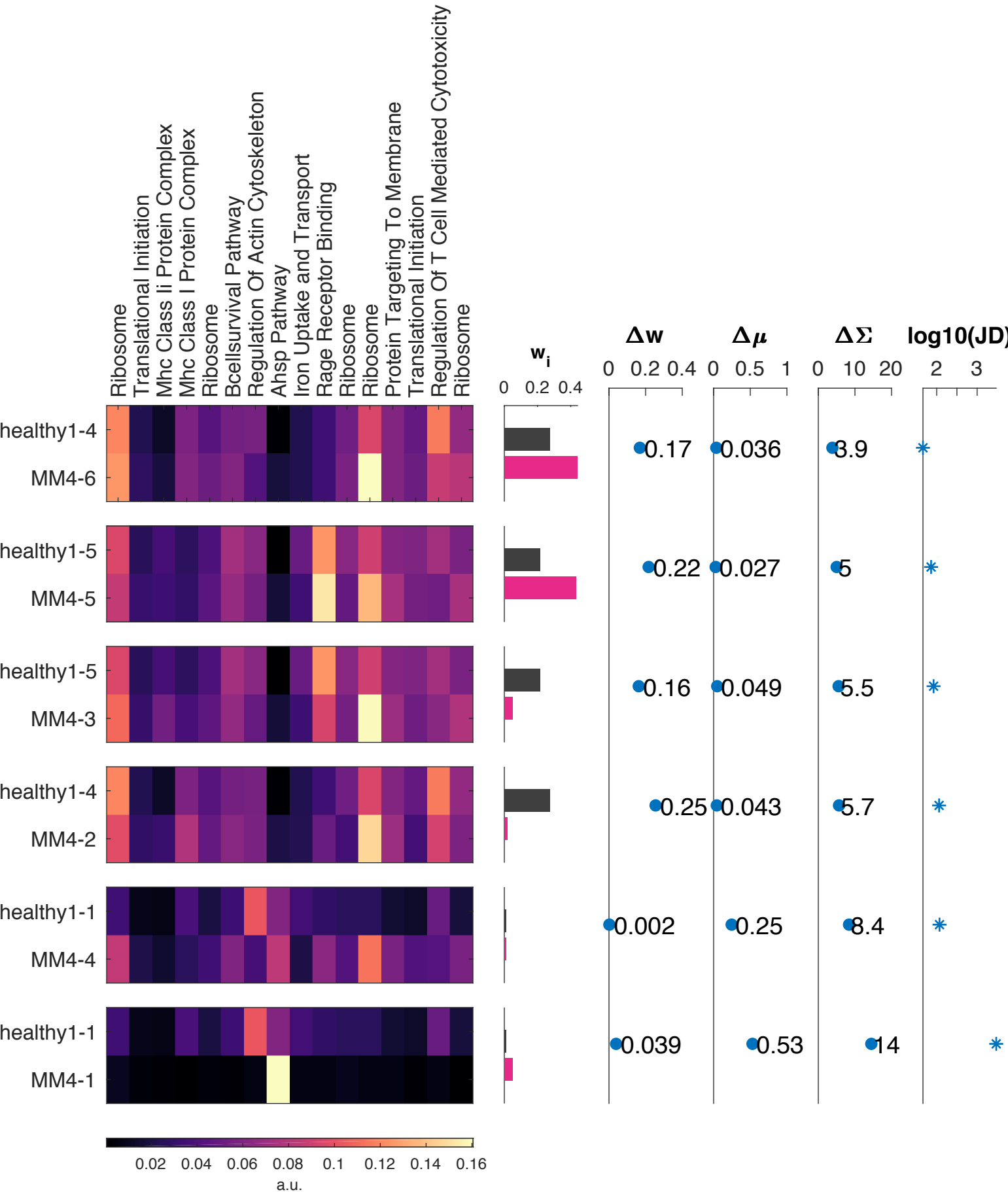
SI Figure 10



SI Figure 11

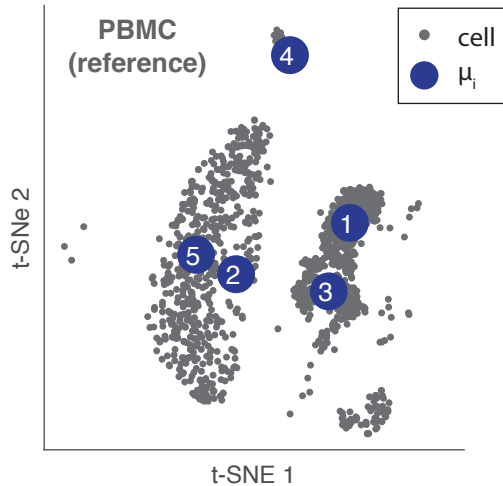


SI Figure 12

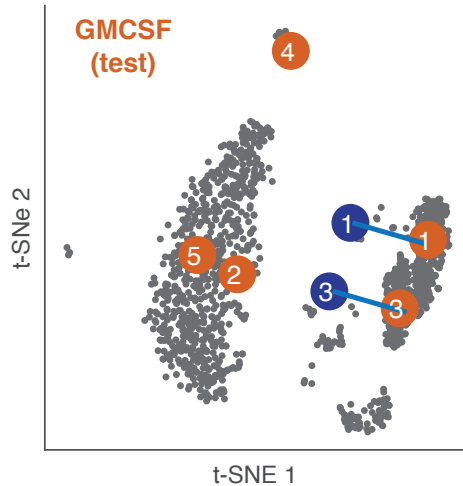


SI Figure 13

a



b



c

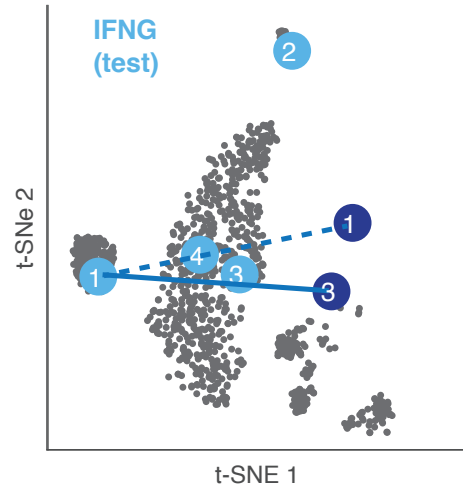


Table 1: Patient Information

Patient Label	Age	Gender	Ethnicity	Disease	Medications
healthy1	58	Male	African American	none	N.D.
healthy2	43	Female	Caucasian	none	N.D.
MM1	51	Male	Caucasian	multiple myeloma	N.D.
MM2	43	Female	African American	multiple myeloma	N.D.
MM3	65	Male	Caucasian	multiple myeloma	Carflizomib/Zometa every 3 months Lisinopril Acyclovir
MM4	67	Male	Hispanic	multiple myeloma	Revlimid and Dexamethasone Hydrocodone Acetaminophen Vitamin D3 Diphenoxylate-atropine Calcium Alendronate Acyclovir Amlodipine Fenofibrate Warfarin